# IJACSA

WHERE WISDOM SHARES

## International Journal of Advanced Computer Science and Applications

Volume 15 Issue 1

January 2024

SAI

www.ijacsa.thesai.org

# Editorial Preface

## From the Desk of Managing Editor...

It may be difficult to imagine that almost half a century ago we used computers far less sophisticated than current home desktop computers to put a man on the moon. In that 50 year span, the field of computer science has exploded.

Computer science has opened new avenues for thought and experimentation. What began as a way to simplify the calculation process has given birth to technology once only imagined by the human mind. The ability to communicate and share ideas even though collaborators are half a world away and exploration of not just the stars above but the internal workings of the human genome are some of the ways that this field has moved at an exponential pace.

At the International Journal of Advanced Computer Science and Applications it is our mission to provide an outlet for quality research. We want to promote universal access and opportunities for the international scientific community to share and disseminate scientific and technical information.

We believe in spreading knowledge of computer science and its applications to all classes of audiences. That is why we deliver up-to-date, authoritative coverage and offer open access of all our articles. Our archives have served as a place to provoke philosophical, theoretical, and empirical ideas from some of the finest minds in the field.

We utilize the talents and experience of editor and reviewers working at Universities and Institutions from around the world. We would like to express our gratitude to all authors, whose research results have been published in our journal, as well as our referees for their in-depth evaluations. Our high standards are maintained through a double blind review process.

We hope that this edition of IJACSA inspires and entices you to submit your own contributions in upcoming issues. Thank you for sharing wisdom.

**Thank you for Sharing Wisdom!**

# Editorial Board

# CONTENTS

(xiii)

# Reliability Evaluation Framework for Centralized Agricultural Internet of Things (Agri-IoT)

Fatoumata Thiam[1], Maïssa Mbaye[2], Maya Flores[3], Alexander Wyglinski[4]

Laboratoire d'Analyse Numérique et Informatique(LANI), University Gaston Berger (UGB), Senegal[1,2]

Department of Electrical and Computer Engineering, Worcester Polytechnic Institute (WPI), USA[3,4]

*Abstract*—This paper presents a holistic reliability evaluation framework for Agri-IoT based on real-world testbed and mathematical modeling of network failure prediction. A testbed has been designed, implemented, and deployed in the real-world in the experimental farm at Saint-Louis/Senegal as a representative area of Sahel conditions. Data collected has been used for real-world reliability analysis and to feed mathematical modeling of network reliability based on energy and environmental conditions data with Kaplan Meier and Nelson Aalen estimators. Key factors affecting the network's lifespan, such as network coverage and density, are explored, along with a comprehensive evaluation of energy consumption to understand node discharge rates impact. The survival analysis, employing Kaplan-Meier and Nelson-Aalen estimators, establishes network stability and the probability of node survival over time. The findings contribute to the understanding of Agri-IoT reliability in a real-world Sahel environment, offering practical insights for system optimization and environmental challenge mitigation in real-world deployments.

*Keywords*—*Energy; IoT; reliability; real-world testbed; optimization; Agri-IoT*

## I. INTRODUCTION

IoT is increasingly adopted in agriculture (Agri-IoT) to improve management methods, performance, and productivity of agricultural farms[1]. In particular, Agri-IoT provides tools for input management, automatic irrigation management, remote control of agricultural fields, yield forecasting, and input prediction.

Deployed Agri-IoT systems are very common in many European and American countries (North and South) and Asia [2], [3]. In developing countries, especially in Africa, more particularly in the Sahelian zone, the deployment of such systems is needed for food security purposes [4]. However, to reproduce the same results of Agri-IoT as elsewhere, reassessments of the reliability of IoT are needed in the specific environment of the Sahel.

Indeed, the Sahel has many white areas without a network, the sector power is not ensured in agricultural areas, and a dry, dusty environment is sometimes rainy which has an impact on the reliability of the electronic devices, the sensors, and the maintenance in the event of a fault is not always possible.

Limited power resources in the IoT network devices, can cause failures due to internal instability or external disturbances. Identifying vulnerabilities and using reliability analysis [5] and fault diagnosis techniques [6] can enhance system stability and resilience.

This reliability assessment is necessary for good network lifespan and data sensor reliability knowing food security issues if an Agri-IoT device fails. More importantly, most of the proposed research is theoretical or in lab tests.

The reliability of IoT systems can be assessed using different metrics such as power consumption and its impact on the network lifetime, network lifespan, sensors data trustworthiness (error rate), network availability, and environmental impact.

This paper aims to provide a holistic reliability analysis framework on a Sahel area based on real-world data collection from the experimental farm and mathematical modeling of network failure prediction.

The main contributions of this paper are: begin

- Design, implement, and deploy real-world Agri-IoT deployment in the experimental farm at Gaston Berger University in Saint-Louis (Senegal).

- Collect real data for reliability evaluation.

- Experiment with different environmental constraints in IoT operation.

- Modeling mathematically the network reliability based on energy and environmental conditions data with Kaplan Meier and Nelson Aalen estimators.

This paper is organized as follows: Section I gives a brief background on IoT and the problem statements. Section II presents a related work on IoT power evaluation and modelization strategies. Section III looked at reliability assessment tools, especially the mathematical techniques on lifetime evaluation. Section IV presents an experimentation of Agri-IoT system deployment and, in that same way, presents the results of our power evaluation and survival analysis for a centralized Agri-IoT Network in real-life deployment. Section V presents the evaluation, analysis, and discussion and Section VI, finally presents the Conclusion.

## II. RELATED WORK

Deploying Agri-IoT systems is challenging in terms of hardware, network architecture density (topology and density), type of power supply, and environment. Battery drain is a parameter of reliability related to network lifespan.

Regarding that, several research [7], [8], [9] evaluated power consumption in different LPWAN technologies. In [7], an analytical approach is proposed to assess individual sensor node power consumption, providing insights for optimizing

sensor node design with a focus on energy autonomy. In [8], the authors presented an energy model for NB-IoT, considering power-saving modes and discontinuous reception mechanisms. In [9], the authors compared the power consumption impact of SF7 and SF12 and their respective applications in a LoRaWAN-based IoT system.

In [10], the authors evaluated the energy usage of LPWAN wireless technologies (LoRaWAN, DASH7, Sigfox, and NB-IoT) to determine battery lifetimes. They found that actual battery lives can differ from ideal scenarios and provide insights on selecting appropriate technologies and battery capacity to improve IoT applications.

Banti et al. [11] identified challenges in designing an energy-efficient LoRaWAN communication protocol. Their findings guide research toward a GreenLoRaWAN protocol that is robust, scalable, and energy efficient. The authors emphasized the need for independent power sources for IoT nodes, as studies often overlook network lifetime constraints by assuming gateways are connected to the grid.

In [12], the authors proposed a node lifetime estimation approach applicable to both static and dynamic loads, investigating the influence of parameters such as self-discharge, discharge rate, age, and temperature on Alkaline and Nickel-Metal-Hydride (NiMH) batteries.

Based on reviews, many energy models simplify analysis for manageability, potentially overlooking real-world complexities and leading to inaccuracies in predictions. Some studies assume grid power for gateways, but in the Sahel, where reliable power is rare, overlooking network lifetime constraints in IoT system design can affect practicality. Dynamic factors impacting energy consumption over time, like environmental conditions and hardware variations, may not align with assumed ideal conditions. In this study, battery consumption estimation relies on node voltage measurement and payload tests, using real-world testbed data under actual conditions in the Sahel. This approach is based on modeling mathematically of network reliability on Kaplan Meier and Nelson Aalen estimators.

### III. BACKGROUND PROBABILISTIC IoT RELIABILITY PREDICTION

In statistics, survival analysis [13], [14], [15] is for analyzing life expectation and lifetime based on an event occurring such as the death or failure of a mechanism in a system. This topic is also called reliability analysis.

Reliability analysis determines the probability of a population that survives past a specific time, the rate of the population that survived or died at any given time, or how an event impacts the population's lifetime. In the Agri-IoT context, it enables modeling network nodes' survival by predicting it.

In mechanical systems, determining the cause of death or failure of a component is required. In survival analysis, this is considered an "event" and involves time-to-event data. Time is defined as the beginning or end of an observation period. Censored observation focuses on an individual's survival time, even if the information is incomplete or imprecise.

In 1958, Kaplan and Meier introduced the Kaplan-Meier estimator [15], which has since become widely used for estimating and summarizing survival curves. This approach is the most common way to estimate and summarize survival curves for being a highly cited statistic paper. The survival function, $S(t)$, gives the probability that a device or a mechanism survives after a specific time $t$. The non-parametric estimator functions will be used to analyze the survival data from field experiments and to evaluate two related probability paradigms.

1) The survival probability using the Kaplan-Meier estimator (KME).
2) The Hazard rate denoted $H(t)$ using the Nelson-Aalon estimator(NAE) [16], [17].

The Kaplan-Meier distribution estimates the probability of the event of interest not happening at time $t$. At the same time, the NAE Hazard probability gives a visualization of the event occurring on the subject within an interval of time.

$S(t)$ is the probability that a given population member has a lifetime greater than $t$. For a sample of size $N$ in this population, the time until each death of the members of population $N$ is represented as follows:

$$t_1 \leq t_2 \leq t_3 \leq \cdots \leq t_N. \tag{1}$$

The KME function $S(t)$ is express as following :

$$\hat{S}(t) = \prod_{t_i < t} \frac{n_i - d_i}{n_i} \tag{2}$$

The NAE function $H(t)$ is expressed as following :

$$\hat{H}(t) = \sum_{t_i < t} \frac{d_i}{n_i} \tag{3}$$

where $n_i$ is the number of subjects "at risk" just before time $t_i$, and $d_i$ is the number of deaths at time $t_i$. The KME is not used to estimate the cumulative hazard, plus the NAE hazard function has a better small-sample performance than the KME [18] function. An empirical comparison of the two solutions has been broadly studied by Colosimo *et al.*[19].

The IoT systems are hardware and software deployed to monitor specific parameters. In some cases, human intervention can be a requirement for provisioning and maintenance reasons or as an oracle. Then, the reliability of IoT systems can be assessed in three ways:

- Hardware reliability;

- Software or operation reliability;

- External reliability.

Hardware reliability is the failure rate of the hardware component system. It has been assessed in many ways in the literature [20], [6], [21]. The software reliability involves protocols, logistic support, and the system's operationality. External reliability is correlated to the maintenance of the system or human intervention and the impact of outside parameters in the system. The system reliability can be expressed as follows:

$$R_{system} = \frac{R_{gateway}(\sum (R_{mn}) + R_{software} + R_{human}}{3})$$

(4)

Where $R_{gateway}$ is the reliability of the gateway, which centralizes the network and all the system depends on. $R_{mn}$ is the reliability of the network participant, the member nodes. $R_{software}$ is the software reliability of participants. $R_{human}$ is the human reliability.

In this study, the $R_{software}$ is not considered, same for the $R_{human}$ on the system reliability. Consequently, in this work, several evaluation and prediction results are achieved. This is based on an IoT system's energy consumption behavior during its operation with on-field collected data from the testbed: Node Lifespan; Network lifespan; Impact of activity level on the system; Network density, and participants' death over time for maintenance and human intervention time prediction.

The next section presents the simulation part and how the reliability analysis is executed and correlated to the on-field implemented system's observations.

## IV. CONTRIBUTIONS: EXPERIMENTATIONS AND MODELING

### A. Context and Test Bed Architecture Design

The test bed site is located at Gaston Berger University's (UGB) experimental farm in northern Senegal, specifically in Saint-Louis (see Fig. 1).

This site has been set up for practical training and research activities focusing on animal and crop production techniques. It covers an area of 26 $ha$ and is irrigated with a drip irrigation system and two submerged pumps in the Djeuss River, Senegal River, with flow rates of $160\ and\ 196\ m^3\ per\ hour$. It has a storage and recovery basin measuring 15 meters by 10 meters by 8 meters, a filtering station with sand tanks, multiple control heads, distribution ramps, and conduits supporting the drip emitters. This infrastructure enables a controlled water supply for agricultural experimentation on the farm.

The Agi-IoT test bed installation was done between August and September 2021. Meteorological conditions during this period are characterized by the rainy season, sunny days, and daytime temperatures ranging from 28°C to 33°C during the day and 24°C and 27°C during the night, while relative humidity levels consistently exceed 70% based on on-field measurements and from National Agency for Civil Aviation and Meteorology [22], [23].

The architecture of the test bed consists of one Gateway and four sensor nodes deployed into an okra exploitation. The system provides sensing and analyzing functions (Fig. 2). The aim is to monitor environmental parameters (temperature and soil humidity) and infer the need for irrigation actioning.

The sensing part consists of a digital temperature and humidity sensor, a controller, and an energy source from a finite battery. Moreover, the sensor used in this work incorporates an integrated temperature sensor to gauge the soil's temperature accurately. The collected data is transmitted to the gateway using BLE technology. Bluetooth Low Energy (BLE)

technology is known for providing reliable communication over relatively short distances.

The IoT's Cisco Reference Model defines the gateway as a Level 3 and 4 component that acts as a hub for receiving all data collected from sensor nodes [24]. The gateway processes the packets to extract relevant information, generates customized irrigation plans for each sensor node, and transmits refined data to a database. This makes it accessible to various applications and services for the end-users.

### B. Hardware Design, Implementation, and Deployment

The IoT Network test bed consists of a gateway and sensor nodes. These components encompass a controlling board, a communication module (integrated into the controller), sensors, and the power supply. The nodes, including gateway and sensor nodes, are implemented using the Raspberry Pi 3B+ as the main electronic board (Fig. 3). Nodes use BLE 4.2 for network, power bank, solar panels, soil sensors for data collection, UM25C Meter for energy monitoring, and an 8-inch display, keyboard, and mouse for RPIs.

The Cypress CYW43455 Bluetooth Chipset and BLE 4.2 are used for the Network Communication interface between the nodes. The BLE 4.2 module is built in the RPi 3B+ for communication among the Network's nodes. In this context, the Gateway is the master, and simple nodes are slaves.

A solar power bank with a battery capacity of $36,000mAh$ is connected to the nodes (gateway and simple node). The power bank has a USB output for the devices' power supply. A $24$ by $6$ cells poly-crystalline silicon solar panel is coupled to the gateway battery. The criteria for choosing the battery model were durability, robustness to the outdoor conditions, and water resistance in accordance with the deployment area conditions.

The UM25C Meter [25] is a device that measures electrical quantities, such as voltage, current, resistance, capacitance, and temperature. It has a large display screen and can be connected to a computer or smartphone via Bluetooth. It is connected to the RPI via USB 2.0. It is used to monitor our field ambient temperature and nodes' electrical metrics quantities.

A WiFi connection was also used to establish Internet connectivity for the gateway. However, this connection was deliberately not kept continuously active in order to reduce power consumption. Instead, the connection was established twice daily to transmit data to the database and was subsequently disconnected. The data acquired from the sensor node was stored using Google Sheets as the storage platform.

In the experimental farm, five(5) nodes have been deployed on the field in a star topology to collect data. The Gateway is in the center with a solar panel and enclosure at $1m$ up to the ground (Fig. 4). It coordinates communication and has two network interfaces: One with BLE to communicate with the sensor nodes in WI-FI to be connected to the Internet. The four(4) Sensor nodes are in the perimeter of a chosen area to collect temperature and humidity and send data to the gateway (Fig. 5).

Fig. 1. Location of the outdoor experimental farm of UGB located at UGB, Saint-Louis, Senegal.



Fig. 2. Agri-IoT deployment architecture diagram.

## V. Evaluation, Analysis, and Discussion

### A. Reliability Evaluation Results

This section presents the main results and learned lessons from the testbed related to reliability. The on-field testbed deployment further highlights the importance of designing and simulating scenarios with detailed realism, ensuring accurate and representative results regarding the hostile Sahel environment.

Table I presents the status of various parameters that impact energy consumption in a node. The gateway plays the role of a server. Several experiments were realized to

TABLE I. The Gateway Internal Activities

|  | Status |
|---|---|
| Server | Gateway |
| Platform RPI | 3B+ |
| Running Tasks | 156 |
| User | Python3.7 |
| Memory | 873.3/1000 |
| Network | BLE |
| Display | None |

observe the system's activity level and running tasks. The running software is implemented in Python with the BLE communication protocol.

Fig. 3. Gateway and sensor basic components lab integration and testing. (1) Raspberry Pi 3B+ with built-in BLE 4.2, (2) UM25C Meter, (3) Battery power bank, (4) STEMMA soil moisture and temperature sensor.



Fig. 4. (1) Gateway's initial deployment: Solar power bank responsible for energy supply. This setup lacks protection against the sun, significantly impacting board operations due to enforced sleep mode during high temperatures within the enclosure. (2) Enclosed gateway setup featuring an RPI, a body with a fan, and a UM256C meter for energy behavior monitoring during node activities.

*1) Network Coverage:* Considering a node as dead once it stops operating. The operation time of the field experiment is 72 hours testing period. The maximum communication distance for the BLE connection was about 11 meters in line of sight. The BLE range of the Raspberry Pi is limited without an external antenna. For any distance beyond 7 to 10 meters or another use case that involves reasonable communication distance, adding an antenna might suit best.

*2) Impact of Density in Network Lifespan:* It is very important to have a global view of the impact of network density on the network to provide energy-saving and optimized operation. Fig. 6(2) shows how the network size impacts the

network's longevity. The experiment was run five times and plotted the mean values of the network lifetime against the network size. Then, observations were made based on the activity and the node's role(sensor node or gateway) to obtain a concise evaluation of its lifespan. Node density was set to 5, 10, 15, 20, and 25 participant nodes with one gateway in a centralized architecture.

Two important facts were deducted from the evaluation of Network Density over Node Lifetime. The red curve represents the lifespan of 'at-risk' nodes over time, which are nodes that are more prone to failure. On the other hand, the blue curve illustrates the lifespan of 'died nodes' over time, which are

Fig. 5. Field deployment of the testbed IoT network, emphasizing a sensor node, its energy source.



Fig. 6. (1) Node discharge rate effect on lifetime. (2) Influence of network density on node lifetime: The red curve represents the lifespan of 'at-risk' nodes over time, while the blue curve illustrates the lifespan of 'died nodes' over time.

the nodes that have experienced failure events. This analysis provides a valuable understanding of the network's stability, reliability, and overall performance.

*3) Impact of energy in network lifespan and reliability :* The reliability evaluation encompassed a comprehensive assessment of discharge rate impact and inclusion of factors such as transmission, reception, idling periods, duty cycling, data, and the node's designated tasks. This comprehensive evaluation allowed us to estimate the node's lifespan, as depicted in Fig. 6(1). Thus, the figure illustrates that with a discharge rate of up to $0.054Ah$, the projected node lifetime could extend beyond $25hours$. As demonstrated through indoor testing, this prolonged node lifespan translates to approximately $22hours$ of sustained gateway operation connected to a finite power source of $36000mA$. This is a way to predict network failure and anticipate self-healing methods to minimize network

failures.

This observation highlights the critical role of gateway efficiency and longevity in overall network performance by deepening our understanding of the correlation between node discharge rates and lifespan.

The battery capacity can be expressed as a function of the time it takes to charge fully (see Eq. 5):

$$C_n = \frac{A(charging\ or\ discharging)}{Capacity} \qquad (5)$$

- Where $C_n$ or $C/n$, expressed in ampere-hours $(Ah)$, measures the speed of charging or discharging the battery, and the $n$ stands for the number of hours the discharging takes.

- $A$ is the electrical current.

- The capacity is the amount of current held in the battery; it is different from the power.

To find out this quantity of energy (which is expressed in Watt-hours - $Wh$), the capacity must be multiplied by the voltage of the battery:

$$Ah \times V = Wh \qquad (6)$$

The current $(A)$ can fluctuate depending on the charging style or type, and different parameters like heat, dust, wear-out, and activity load can affect battery charging and discharge speed. Hence, All battery parameters are affected by the battery charging and recharging cycle. The current can be expressed in Eq. 7.

$$A = \frac{capacity}{H} \qquad (7)$$

From Eq. 5 and 7,the lifetime is deduce as follows:

$$H = \frac{1}{Cn} \qquad (8)$$

*4) Mitigating environmental-related issues::* Nodes went to sleep after sending measurements, making data collection impossible. The gateway received data only the next morning. No protection against the sun heat in Fig. 3 affected performance. High temperatures triggered sleep mode, posing a deployment challenge.

The temperature was $46°$C in the field. High temperature triggers sleep mode and reduces charging efficiency. More energy is required for proper recharging. High temperatures slow down charging. Our batteries come with a less-efficient amorphous solar panel.

The power bank's operational temperature range spans from $0°$ to $45°$, with a cautionary recommendation not to exceed $60°$. During the deployment phase, high temperature consequently impacted the battery's performance. This explains the occurrence of node failures in the field due to the influence of heat.

Furthermore, an important observation concerning the new gateway's design. It featured a wider solar panel (refer to Fig. 4(1)) that served as a protection against sun heat, effectively preventing node overheating. The system's efficiency is improved by adjusting the battery supply and raising the solar panel. This created better airflow, resulting in a cooling effect. The change was needed because inconsistent data reception was experienced from the end nodes. Initially, attempts were made to enhance the gateway's power supply, but upon further investigation, it was determined to be unnecessary. As a result, the gateway functions optimally in cooler environments with average temperatures ranging from $25°$ to $30°$.

## B. Network Survival Analysis

A thorough survival analysis was performed to explore the correlation between network density and the lifespan of the network. This investigation enabled us to make accurate predictions about the network's behavior in response to a 10% increase or decrease in workload or network size. Fig. 7 is a visual representation of our findings. Network expansion impacts payload and data processing at the gateway. This data was analyzed using Python 3.7 and the Lifelines library, with the Kaplan Meier distribution fitter method - a widely accepted approach.

During the investigation, time-to-event data was analyzed using two non-parametric survival function estimator techniques: the Kaplan-Meier Estimator (KME) and the Nelson-Aalen Estimator (NAE) to consecutively estimate $S(t)$ and $H(t)$.

The KME was used to estimate the probability of nodes surviving over time. This estimator is represented as a step function, which shows discontinuities at the occurrence of events. Initially, the probability of a node's survival remains at 100%, indicating network stability during the first 18 hours, between the intervals of $t_0$ and $t_{18}$, as shown in Fig. 8(2). This interval represents a period of network stability. Additionally, the median survival time indicates that nodes can operate reliably for at least 22 hours, providing valuable information for proactive network maintenance planning.

Due to the limited amount of data, the Kaplan-Meier estimate was not applicable. Consequently, the Nelson-Aalen estimator, more accurate for a limited amount of data, was utilized to assess the cumulative hazard rate. This involves calculating the total number of node failures during specific time intervals to determine the cumulative count. Refer to Fig. 8(2) for the results.

Our analysis found no "early node mortality." Our metric for assessing node performance was battery depletion. Failures became more prevalent as the network aged, suggesting that nodes performed well initially but gradually became less reliable over time.

In the scope of the study, the metric related to battery depletion in the network component was closely monitored. It is important to note that no subjects were censored in this investigation. This can be visualized in Fig. 8(1) and 8(2).

## VI. CONCLUSION

This paper evaluates the reliability of an Agri-IoT system deployed in the challenging Sahel environment. Accurately simulated on-field scenarios provided valuable insights. Monitoring the gateway's internal activities revealed critical data impacting node energy consumption. Network coverage and density emerged as key factors affecting the network's lifespan. The comprehensive evaluation of energy consumption provided crucial insights into node discharge rates and their correlation with lifespan, aiding in predicting and mitigating network failures.

The study addressed environmental challenges, highlighting the impact of high temperatures on node performance and the successful design adaptation of the gateway for improved efficiency in cooler environments.

Fig. 7. Network density and duty-cycling impact on the network lifetime.



(1)



(2)

Fig. 8. (1) Kaplan Meier estimate survival function, $S(t)$ described in Eq. 2. It is the probability that a node survives from the time it is switched on to a specific future time $t$. (2) Nelson-Aalen Estimator cumulative hazard rate. It is denoted $H(t)$ described in Eq. 3. It represents the probability a node, in our case, who is under observation at a time $t$, has died at that time.

Survival analysis, utilizing Kaplan-Meier and Nelson-Aalen estimators, established network stability and the probability of node survival over time. The absence of "early node mortality" indicated initial reliability, gradually decreasing over time.

These results enabled to better understanding of Agri-IoT reliability in typical Sahel Environment, providing practical insights for optimizing system performance and addressing environmental challenges in real-world deployments.

### REFERENCES

[1] S. Rudrakar and P. Rughani, "Iot based agriculture (ag-iot): A detailed study on architecture, security and forensics,"

*Information Processing in Agriculture*, Sep. 2023. [Online]. Available: https://doi.org/10.1016/j.inpa.2023.09.002

[2] V. Saiz-Rubio and F. Rovira-Más, "From smart farming towards agriculture 5.0: A review on crop data management," *Agronomy*, vol. 10, no. 2, p. 207, 2020, accessed November 25, 2023. [Online]. Available: https://www.mdpi.com/2073-4395/10/2/207

[3] D. Pivoto, P. D. Waquil, E. Talamini, C. P. S. Finocchio, V. F. Dalla Corte, and G. d. V. Mores, "Scientific development of smart farming technologies and their application in brazil," *Information Processing in Agriculture*, vol. 5, no. 1, pp. 21–32, March 1 2018. [Online]. Available: https://doi.org/10.1016/j.inpa.2017.12.002

[4] R. K. Goel, C. S. Yadav, S. Vishnoi, and R. Rastogi, "Smart agriculture-urgent need of the day in developing countries," *Sustainable Computing: Informatics and Systems*, vol. 30, p. 100512, June 1 2021. [Online]. Available: https://doi.org/10.1016/j.suscom.2021.100512

[5] A. Azhdari, M. A. Ardakan, and M. Najafi, "An approach for reliability optimization of a multi-state centralized network," *Reliability Engineering and System Safety*, vol. 239, p. 109481, 2023.

[6] Z. Gao, C. Cecati, and S. X. Ding, "A survey of fault diagnosis and fault-tolerant techniques—part i: Fault diagnosis with model-based and signal-based approaches," *IEEE Transactions on industrial electronics*, vol. 62, no. 6, pp. 3757–3767, 2015.

[7] H. Rajab, H. Al-Amaireh, T. Bouguera, and T. Cinkler, "Evaluation of energy consumption of lpwan technologies," *EURASIP Journal on Wireless Communications and Networking*, vol. 2023, no. 1, p. 118, 2023.

[8] A. K. Sultania, P. Zand, C. Blondia, and J. Famaey, "Energy Modeling and Evaluation of NB-IoT with PSM and eDRX," in *2018 IEEE Globecom Workshops, GC Wkshps 2018 - Proceedings*. Institute of Electrical and Electronics Engineers Inc., feb 2019.

[9] T. G. Durand, L. Visagie, and M. J. Booysen, "Evaluation of next-generation low-power communication technology to replace GSM in IoT-applications," *IET Communications*, vol. 13, no. 16, pp. 2533–2540, 2019. [Online]. Available: www.ietdl.org

[10] R. K. Singh, P. P. Puluckul, R. Berkvens, and M. Weyn, "Energy consumption analysis of lpwan technologies and lifetime estimation for iot application," *Sensors*, vol. 20, no. 17, p. 4794, Aug 2020. [Online]. Available: http://dx.doi.org/10.3390/s20174794

[11] K. Banti, I. Karampelia, T. Dimakis, A.-A. A. Boulogeorgos, T. Kyriakidis, and M. Louta, "Lorawan communication protocols: A comprehensive survey under an energy efficiency perspective,"

*Telecom*, vol. 3, no. 2, p. 322–357, May 2022. [Online]. Available: http://dx.doi.org/10.3390/telecom3020018

[12] W. Rukpakavong, L. Guan, and I. Phillips, "Dynamic node lifetime estimation for wireless sensor Networks," *IEEE Sensors Journal*, vol. 14, no. 5, pp. 1370–1379, may 2014.

[13] G. Rodrıguez, "Non-parametric estimation in survival models," *cited on*, p. 20, 2005, accessed on November 25, 2023. [Online]. Available: https://grodri.github.io/survival/NonParametricSurvival.pdf

[14] L. J. Bain, "Analysis for the linear failure-rate life-testing distribution," *Technometrics*, vol. 16, no. 4, pp. 551–559, 1974.

[15] E. L. Kaplan and P. Meier, "Nonparametric estimation from incomplete samples," *Journal of the American Statistical Association*, vol. 53, no. 282, pp. 457–481, 1958. [Online]. Available: http://www.jstor.org/stable/2281868

[16] O. Aalen, "Nonparametric inference for a family of counting processes," *The Annals of Statistics*, pp. 701–726, 1978.

[17] W. Nelson, "Theory and applications of hazard plotting for censored failure data," *Technometrics*, vol. 14, no. 4, pp. 945–966, 1972.

[18] J. KLEIN, "Small sample moments of some estimators of the variance of the Kaplan-Meier and Nelson-Aalen estimators," *Scandinavian Journal of Statistics*, vol. 18, no. 4, pp. 333–340, 1991.

[19] E. Colosimo, F. v. Ferreira, M. Oliveira, and C. Sousa, "Empirical comparisons between kaplan-meier and nelson-aalen survival function estimators," *Journal of Statistical Computation and Simulation*, vol. 72, no. 4, pp. 299–308, 2002.

[20] J.-J. Lee, B. Krishnamachari, and C.-C. J. Kuo, "Aging analysis in large-scale wireless sensor networks," *Ad Hoc Networks*, vol. 6, no. 7, pp. 1117–1133, 2008.

[21] I. Kabashkin and J. Kundler, "Reliability of sensor nodes in wireless sensor networks of cyber-physical systems," *Procedia Computer Science*, vol. 104, pp. 380–384, 2017.

[22] ANACIM, "Anacim — agence nationale de l'aviation civile et de la météorologie," 2023. [Online]. Available: https://www.anacim.sn/

[23] A. . A. N. de l'Aviation Civile et de la Météorologie. (2021) Bulletin saisonnier jas 2021. [Online]. Available: https://bit.ly/3MZdsSZ

[24] C. Systems, "Iot: A cisco model," *LearnIoT.com*, 2023. [Online]. Available: https://learniot.com/cisco-model

[25] Joy-IT, "Um25c multimeter," 2023. [Online]. Available: https://joy-it.net/en/products/JT-UM25C

# A Hybrid Approach for Automatic Question Generation from Program Codes

Jawad Alshboul, Erika Baksa-Varga

University of Miskolc, Faculty of Mechanical Engineering and Informatics, Miskolc, Hungary

*Abstract*—Generating questions is one of the most challenging tasks in the natural language processing discipline. With the significant emergence of electronic educational platforms like e-learning systems and the large scalability achieved with e-learning, there is an increased urge to generate intelligent and deliberate questions to measure students' understanding. Many works have been done in this field with different techniques; however, most approaches work on extracting questions from text. This research aims to build a model that can conceptualize and generate questions on Python programming language from program codes. Different models are proposed by inserting text and generating questions; however, the challenge is understanding the concepts in the code snippets and linking them to the lessons so that the model can generate relevant and reasonable questions for students. Therefore, the standards applied to measure the results are the code complexity and question validity regarding the questions. The method used to achieve this goal combines the QuestionGenAi framework and ontology based on semantic code conversion. The results produced are questions based on the code snippets provided. The evaluation criteria were code complexity, question validity, and question context. This work has great potential improvement to the e-learning platforms to improve the overall experience for both learners and instructors.

*Keywords*—*Question generation; e-learning; python question generator; semantic code conversion*

## I. INTRODUCTION

Automating question generation has become significant with the increasing trend of online learning and its scalability in recent years. Technical courses like learning programming languages are more popular, and there is a massive demand for such subjects. Questions are the primary approach used to evaluate student knowledge [1]. Therefore, creating questions becomes more challenging as the constant growth of e-learning continues, more courses are created, and the pressure on teachers is high. Intelligent and deliberate questions can enhance student understanding and reduce the gap between theory and practice in programming subjects [2]. For example, the article in [3] monitors the performance and behavior of students who engage in courses with self-assessment methods in programming and problem-solving. The research in [4] observes the decentralized practice by monitoring the intensity and timing of the impact on student learning and problem-solving in programming languages. The research paper [5] addresses interactivity while solving problems in programming languages based on learning objects. The article in [6] tries to enhance the use of digital resources for students and instructors. The research papers in [7] and [8] address the learning objects that can be used in different contexts using

web3. Finally, the article in [9] suggests collaborative learning to help instructors engage students in generating and evaluating questions. The proposed method focuses on translating Python code into text and uses an AI-based framework to generate questions from the text. We also use ontology to connect and conceptualize the logic of the programming language. Applying ontology ensures interoperability with other systems and reduces the overhead on educational platforms. This work contributes to the e-learning platforms and improves the overall experience for instructors of programming languages. It also enhances the learning path for students who like to learn and do exercises without repeating the same questions. The outcome of this research is to generate meaningful questions based on Python code to assist instructors in creating more questions in a timely manner, thus ensuring students proper learning of the potential programming language. Unlike similar works, most recent research focuses on generating questions from text, while some research focuses on generating questions from visuals or images [10]. This work focuses on generating questions from code snippets using semantic relations to extract the concepts. Generating questions from unconventional sources, such as code snippets, becomes important in providing a better learning experience to large groups of students, especially when dealing with limited information.

### A. Research Goal

The main goal of this research is to assist instructors and students in properly evaluating student performance by generating Python-based programming questions from existing materials (i.e., code snippets). The automatic question generation from code snippets will add the possibility of generating a different set of questions based on the same code snippet. Therefore, it leads to a better understanding of the given topic.

### B. Research Objectives

To achieve the primary goal of this research, the following objectives are needed:

*1)* Implement a framework that can interpret Python programming language into text.

*2)* Enable the framework to comprehend the text and build connections between the programming structures and the semantic concepts.

The rest of the paper is organized as follows: Section II describes related work and some existing approaches. Section III details the question generation framework implemented in this research, and Section IV shows the results. Section V

presents a discussion that summarizes the results. Finally, Section VI concludes the paper and mentions future work.

## II. RELATED WORK

The question generation process is a relatively complex and challenging task. It requires adequate experience, high knowledge of the material, and time. With the emergence of online learning, it has become a necessity. The first types of question generation models, such as syntax-based, semantic-based, and other models, started in 2014 [11].

Ontologies are a powerful tool for standardizing knowledge representation, which can be helpful in a wide range of domains, including e-learning. By modeling learning materials with ontologies, it is possible to create more personalized and effective learning experiences, allowing learners to achieve their goals more efficiently [12].

Domain knowledge models can be extremely useful in representing knowledge in a standardized and structured manner, aiding the teaching and learning processes. Python and Owlready2 are used to create the model implementation. Python is a popular programming language for various applications, including machine learning and data analysis. It includes many libraries and frameworks for developing sophisticated software systems. Owlready2 is a Python-based ontology library that simplifies creating and manipulating ontologies in Python code. The researchers created a flexible model that can be used to represent knowledge in a way that can be easily integrated into e-learning systems by implementing the domain knowledge model using Python and Owlready2. It could help develop adaptive learning systems that can tailor the learning experience to the needs of individual students [13].

Despite its importance, implementing question-generation approaches for programming concepts is partially applied in the modern world. Programming languages are an essential topic in computer science and software development, and there is a great demand for effective and efficient ways to teach programming concepts. Developing question-generation approaches for programming languages makes it possible to create many practice questions that can be used to reinforce learning and test understanding [14].

To aid students in their learning, Urazova [15] discusses the development of a system to automatically generate questions regarding UML database design and evaluate student responses. The system generates questions and evaluates student responses using artificial intelligence and methods for natural language processing. The goal is to give students a valuable and practical tool to assess their knowledge of and develop their expertise in UML database design.

A study by Russell in [16] investigated the application of automated code-tracing exercises for teaching introductory programming (CS1) courses. Code tracing is a teaching method in which students mentally run code, follow the control flow, and note the values of variables at each stage. The researchers used an automated system to create code-tracing tasks and assess student responses. The purpose of this system was to serve as a tool for teaching students programming ideas and problem-solving techniques. The researchers evaluated the system's efficacy through studies and student polls and contrasted it with more conventional teaching techniques. The article details the possible advantages of automated code-tracing exercises in CS1 courses and the challenges and limitations that must be addressed.

The use of large language models to automatically provide programming tasks and related code explanations has been used recently. A solution using artificial intelligence was developed to help teachers and instructors construct and deliver efficient programming assignments. A sizable language model was trained on text and computer code to provide workouts and explanations for programming. The model was assessed based on the usefulness and quality of the activities and explanations produced through tests and questionnaires. The promise of leveraging extensive language models for automated programming exercise production is discussed in the study, along with the difficulties and restrictions that must be overcome [17].

Automated programming exercise creation and code explanation have several drawbacks, including bias potential, dependency on large language models, limited capacity to assess student comprehension, significant computing needs, and difficulty in generating high-quality training data. The quality of the language models determines how well the exercises and explanations are created, and there is a chance that the models might be biased. The systems automatically assessing student knowledge could not be reliable and might need a lot of computer power. Producing high-quality training data for large language models is challenging and time-consuming. When utilizing these technologies in educational contexts, these constraints must be taken into account [18].

## III. QUESTION GENERATION FRAMEWORK

Question generation involves computer understanding of the available materials to propose plausible questions to students. However, two approaches are usually effective: AI-based or semantic-based [12]. The current work uses a combination of both semantic and AI methods to properly generate questions from code snippets based on semantic code conversion. The primary motivation for using the semantic approach is maintaining concept relations in the programming language keywords to increase system intelligence on the programming language rules. Other approaches would not accurately represent the programming language rules, keywords, and concepts. This section will detail the framework architecture, the technology used, and the approach to generate questions.

### A. Architecture

To generate questions from existing Python code snippets, an interpreter is needed to dissolve the codes into more understandable concepts. Python or any other programming language is constructed using operators, variables, and functions. Operators such as +,-, AND usually do the actual computing. At the same time, variables are used to store values and recall them with operators to perform specific tasks. Functions contain a list of variables, loops, and operators to be executed in order. The ontology will categorize and conceptualize the list of commands (i.e., variables, operators,

etc.) and the relationships between the concepts in the script. It will build an explained version of the code by processing the code line by line and creating semantic relationships based on the input. Subsequently, the translated code is generated and inserted into an AI question generator called "QuestGen" [19]. This model will generate open-ended questions. Fig. 1 shows the framework data flow and its components.

Awareness of existing technologies and software is essential to construct any framework or software. Such awareness can improve productivity and help address many issues that take a long time. As a result, in this work, we implemented a framework using various third-party software. Table I describes this case's environment settings, tools, and applied libraries. As mentioned earlier, we have used the QuestGen AI model, an open-source NLP library dedicated to creating simple question-generation methods. It is on a mission to become the world's most sophisticated question-generation AI by utilizing cutting-edge transformer models like T5, BERT, and OpenAI GPT-2, among others. The primary objective of QuestGen AI is to simplify the question-generation process, providing support to educators, content creators, and learners in developing educational materials. This tool significantly enhances the efficiency of teaching and learning resource development through automation, ultimately facilitating a more effective educational experience.

Before generating questions, the QuestGen AI model expects a text as input. The ontology mentioned next is responsible for converting the snippet code from the Python programming language into text that humans can understand. Subsequently, this model can generate questions based on the inserted text. The software supports four types of questions, and they are as follows:

- Questions with Several Choices (MCQs)
- Boolean (Yes/No) Questions
- Open-ended Questions
- Question Paraphrase

For the current study, only open-ended questions are considered. Since learning a programming language focuses on understanding the content of a code, it is more suitable to use open-ended questions to assess student knowledge properly.

### B. Ontology Design

The ontology is built and compiled using the OWLReady2 library in Python. Such a library would support automating manual activities like adding instances to the ontology. However, the main components and the relationships between concepts should be implemented manually to maintain logical correctness. Translating code into text starts with assigning keywords to ontology classes and describing these keywords. For example, the "=" sign is described in the ontology as an "equal sign", a value of the Assignment subclass in the operator class. The output of the ontology implemented in Python and OwlReady2 is then imported into Protégé for visualization purposes since the visualization is not yet supported on OwlReady2. Fig. 2 shows the ontology design visualization in Protégé.



Fig. 1. Proposed framework architecture.

TABLE I. RESEARCH QUESTIONS AND CORRESPONDING RESEARCH OBJECTIVES

| Name | Description |
|---|---|
| OwlReady2 | Python library to implement Ontology V 0.37 |
| Protege | Software Application for viewing and modifying ontology |
| Jupyter Notebook | IDE to develop the framework |
| QuestGen | AI-based application to generate questions from the text |
| Python | V 3.11.1 |



Fig. 2. Ontology design visualization using protégé.

Logical correctness would enforce semantic meaning on the written script. For example, an "elif" statement syntax is valid in Python. However, it cannot exist without having an "if" statement before it. An "elif" should only be coming after an "if". Furthermore, logical correctness would connect all the keywords and describe the semantic relationship between steps. Most essential aspects of Python programming language in the designed ontology are classified as classes and subclasses. For example, in this study, we have categorized the Python language elements and constructs into four main

classes: Control Structure, Function, Library, and Operator. Each subclass of the Operator class contains several instances that would map each instance to the operator class. Such mapping would assist in enforcing the logical correctness of the translated snippet. Fig. 3 shows an instance definition from the constructed ontology. The ontology's capabilities aim to structure the Python programming language to ensure that the computer can collect vocabulary text about the keywords and build sentences based on the combination of the programming language keywords, which can be fed later into the question generation model. The main limitation is that the ontology should be built manually by adding the explanation of all instances, which can be challenging to implement. Further research is needed to improve this approach. Fig. 4 shows a part of the ontology in Python script.

```
<owl:Class rdf:about="#Subtraction">
  <rdfs:subClassOf>
    <owl:Restriction>
      <owl:onProperty rdf:resource="#has_example"/>
      <owl:hasValue rdf:datatype="http://www.w3.org/2001/XMLSchema#string">Example usage of Subtraction</owl:hasValue>
    </owl:Restriction>
  </rdfs:subClassOf>
  <rdfs:subClassOf>
    <owl:Restriction>
      <owl:onProperty rdf:resource="#has_description"/>
      <owl:hasValue rdf:datatype="http://www.w3.org/2001/XMLSchema#string">Description of Subtraction</owl:hasValue>
    </owl:Restriction>
  </rdfs:subClassOf>
  <rdfs:subClassOf rdf:resource="#Arithmetic"/>
</owl:Class>
```

Fig. 3.   Instance definition of subtraction.

```
# Define subsubclasses data structure
subsubclasses = {
    'Arithmetic': ['Addition', 'Subtraction', 'Multiplication', 'Division', 'Modulus', 'Exponentiation', 'Floor division'],
    'Assignment': ['Equal To', 'Add and Assign', 'Subtract and Assign', 'Multiply and Assign', 'Divide and Assign', 'Modulus ...
    'Comparison': ['Equal', 'Not Equal', 'Less Than', 'Greater Than', 'Less Than or Equal', 'Greater Than or Equal'],
    'Logical': ['and', 'or', 'not'],
    'Identity': ['is', 'is not'],
    'Membership': ['in', 'not in'],
    'Bitwise': ['AND', 'OR', 'XOR', 'NOT', 'Zero fill left shift', 'Signed right shift']
}

# Create operator subclasses and subsubclasses in the new ontology
for operator_type, operator_list in subsubclasses.items():
    with onto:
        operator_type_class = types.new_class(operator_type, (Operator,))
        for operator_name in operator_list:
            operator_name_class = types.new_class(operator_name.replace(" ", "_"), (operator_type_class,))

# ############################################### FUNCTIONS ###############################################
function_sub_classes = ['Built_in', 'Custom_defined']

# Add subclasses for Function
for row in function_sub_classes:
    with onto:
        class_name = row.replace("-", "_")
        NewClass = types.new_class(class_name, (Function,))

# #####################################################CONTROL STRUCTURE###############################################
control_structure_sub_classes = ['Sequence', 'Selection', 'Loop', 'Exception_handling']
control_structure_selection = ['If', 'Else', 'Elif']
control_structure_loop = ['For', 'While', 'Try']
control_structure_exception_handling = ['Except', 'Finally', 'With']

# Add subclasses for Control_Structure
for row in control_structure_sub_classes:
    with onto:
        class_name = row.replace("-", "_")
        NewClass = types.new_class(class_name, (Control_Structure,))
```

Fig. 4.   Ontology in python.

## C. Parser

The parser's job is to detach a block of codes into pieces that can match the ontology based on keywords and custom conditions. These conditions are adjusted depending on the inserted snippets. This model uses the ontology to create sentences. It analyzes keywords in the parser and generates sentences explaining the code. For example, a=10, the parser would create "a is a variable. a value is 10". Fig. 5 illustrates how the code parser algorithm works in the implemented system.

1. **Load the ontology**
2. **Parse Python code and collect keyword explanations**
3. **Split the code into tokens**
4. **Check if a token is a keyword**
5. **Get the explanation of a keyword from the ontology**
6. **Return collected explanations**

Fig. 5.   Algorithm steps of the code's parser explanation.

This parser helps turn Python code (and maybe other types later) into sentences using a set of rules. It maintains whatever logic the ontology possesses about the code. Then, it is fed into the AI model to generate proper questions based on the code interpretation by ontology. The limitation of this parser is that it might struggle with complicated code because it needs specific filters to understand the context and collect the keywords. Fig. 5 describes the steps involved in processing the input and generating results. Initially, the ontology file must be loaded into the environment using an OWL file. Subsequently, the Python source code is provided to the application, where filters extract keywords and retrieve explanations from the defined ontology. Finally, the 'explained code' is passed to the QuestionGenAI framework to generate questions.

## D. Question Generation

Over time, there is a growing demand for question generation, a trend that could significantly alleviate the burden on educators and trainers. This is particularly beneficial for scalable learning formats such as online courses. Many models exist for generating questions from regular text; however, understanding code and generating questions from code snippets is not applied due to its complexity. Code-to-text conversion is a challenging task. However, the semantic relationships between the concepts in the ontology are an excellent solution. Fig. 6 shows the whole procedure to translate code into text. In Fig. 6, the code undergoes validation by a parser checker responsible for scrutinizing its syntax. Once the code is confirmed as error-free, the checker directs it to the ontological translator, acting as the parser within our architecture. This parser transforms the code into coherent sentences, forwarding them to the Question Generator AI model to generate reasonable questions. An explanation of the Question Generator AI model is provided in the subsequent section.

Fig. 6.   Question-generation process.

## E. QuestionGen AI

The QuestGen AI model is an AI model that can generate questions using AI. The QuestGen project is available in an

open-source format [18]. The model is already trained and can generate high-quality questions based on text fed into the model. Instructors can choose the type of question that can be generated; however, we have only applied open-ended questions. The results summarized in the subsequent section show that the AI model can generate reasonable questions based on the input text and its level of clarity.

- Input: The model can process various types of input, including structured, unstructured, and context-based content such as passages, documents, and articles.

- Field of application: The model is tailored to support the education field across diverse disciplines such as science, history, language arts, and more. However, it does not have the capability to execute or generate programming language code.

- Generation method: It is a semantic-based model designed to comprehend inserted text by leveraging concepts and contextual awareness. This procedure is divided into two main steps. Firstly, it begins with entity recognition, wherein the model extracts crucial information such as dates, names, and relationships, employing part-of-speech tagging. Next, the model applies question templates to the extracted information to match the most suitable predefined question template. To improve question quality, various methods are employed, including probabilistic approaches to refine wording and phrasing within the questions.

- Question format: The model can propose various formats, including open-ended, multiple choice, true/false, and short answer.

- Response format: The responses are generated in both text and JSON formats. Each type of question has its own format. For instance, multiple-choice questions prompt the system to produce the question stem and its corresponding answer choices. This distinction applies to all question types, and the resulting output is tailored accordingly.

- Example: The sentence inserted into the model is "The Industrial Revolution was a period of significant economic and technological change that began in Britain during the late 18th century. It marked the transition from agrarian economies to industrialized ones, with advancements in machinery, transportation, and manufacturing."

- The generated questions for a true/false type of question are:

  o 'Is the industrial revolution the same as the 'revolution?',

  o 'Was the industrial revolution a period of change?',

  o 'Was the industrial revolution a revolution in the 18th century?'

## IV. RESULTS

The results are generated in two versions, one utilizing our proposed model and the other without its use (i.e., by directly inserting the code into the QuestGen AI), as depicted in Fig. 7. The implemented framework facilitates the question-generation process, empowering teachers to automatically generate Python programming language assessment questions for testing students' knowledge. Fig. 8 depicts a straightforward code snippet featuring variable definitions. This figure illustrates specific variables alongside their assigned values, incorporated as a script within the ontology. A Python parser is employed to validate the text as proper code before generating any flawed or erroneous questions to mitigate the potential for incorrect syntax within the inserted code. Fig. 9 displays the translated text derived from the code, providing a textual interpretation for each line. The interpreter presents the variable type and specifies the assigned value for each variable. Fig. 10 showcases the outcomes resulting from inserting the aforementioned text into the QuestGen AI model. Fig. 11 can be seen without having a context. The question generator failed to produce any meaningful questions except for the list variable, where it managed to generate a relevant question. However, the AI model could not comprehend all the lines, hence the presence of the ZERO {} symbol. Fig. 12 exhibits a Python code comprising class and object definitions presented as a string and passed through an ontology to translate it into text. Subsequently, this text is fed into the QuestionGen model to generate questions. In the subsequent examples, only the generated questions and context from QuestGen AI will be showcased, omitting the complete outputs. Moving on to Fig. 13, it explains the preceding code snippet depicted in Fig. 12 using natural language, preparing it for input into the AI generator. Following this, Fig. 14 displays the questions generated from the snippet description, demonstrating the relevance of the generated questions. However, Fig. 15 illustrates the outcome of generating questions without providing a snippet description, resulting in improper questions marked by ZERO{} symbols and inaccuracies. This indicates the necessity of providing a description for accurate question generation. In the third example, depicted in Fig. 16, a function is defined to compute the area of a circle based on its radius. This code incorporates arithmetic operations and utilizes Python's 'math' module. Subsequently, Fig. 17 exhibits the output resulting from describing the aforementioned code to input into the AI model. Meanwhile, Fig. 18 displays the generated questions derived from the description of the code snippet involving mathematical operations. Conversely, Fig. 19 showcases a question generated without describing the snippet. The results depicted in all figures are formatted in JSON, containing both the question and its solution. For open-ended questions, the QuestGen model provides the answer alongside the question, excluding the options. It is worth noting that there are warnings due to deprecated libraries utilized by the QuestionGen model, prompting necessary updates by the authors.

Fig. 7.  Generating questions directly from code.



Fig. 8.  A code snippet with variable definitions.

```
xfoo is a string variable and its value is 'foo'
ab is a list variable and it has 2 items
cd is a list variable and it has 3 items
ef is an integer variable and its value is 10
```

Fig. 9.  Generated text from a code snippet.



Fig. 10.  Generated questions for variable definitions.



Fig. 11.  Generated questions without using the proposed approach.



Fig. 12.  Python code for defining classes and objects.

```
Person is a class definition
  __init__  is a method
    name is an instance of the property
    age is an instance of the property
Student is a class definition
  __init__  is a method
    school is an instance of the property
  Student inherits from the Person class
var1 is an instance of the Person class with name 'Jane' and age 25
var2 is an instance of the Student class with name 'John', age 20, and school 'ABC School'
Student inherits from Person
```

Fig. 13.  Generated explanation of the code in Fig. 12.

[{'Question': 'What is person?', 'context': 'Person is a class definition'}]

[{'Question': 'What is __init__?','context': '__init__ is a method'}]}

[{'Question': 'Name is an instance of a property?','context': 'name is an instance of the property'}]

[{'Question': 'What is age an instance of?','context': 'age is an instance of the property'}]

[{'Question': 'What a student a class definition?','context': 'Student is a class definition'}]

[{'Question': 'What is __init__?','context': '__init__ is a method'}]

[{'Question': 'What is a school?','context': 'school is an instance of the property'}]

[{'Question': 'What class does a student inherit from?','context': 'Student inherits from the Person class'}]

[{'Question': 'What is var1 an instance of?','context': "var1 is an instance of the Person class with name 'Jane' and age 25"}]

[{'Question': "What is the instance of the Student class with name 'John', age 20, and school 'ABC School'?",'context': "var2 is an instance of the Student class with name 'John', age 20, and school 'ABC School'"}]

[{'Question': 'Who does a student inherit from?','context': 'Student inherits from Person'}]

Fig. 14.  Generated questions for the more advanced snippet.

ZERO{}
[{'Question': 'What is the age of the person in def __init__?','context': 'def __init__(self, name, age):'}]
ZERO{}
[{'Question': 'What does age mean?',  'context': 'self.age = age self.age = age'}]
ZERO{}
[{'Question': 'What is the age of the child?','context': 'def __init__(self, name, age, school):'}]
[{'Question': 'What is super().__init__(name, age)?','context': 'super().__init__(name, age)'}]
[{'Question': 'What is self.school?','context': 'self.school = school self.school = school'}]
ZERO{}
[{'Question': 'What is the value of student(John, 20, "ABC School")?','context': 'var2 = Student("John", 20, "ABC School")'}]

Fig. 15.  Generated questions without using the proposed model.



Fig. 16.  Code snippet containing a function and arithmetic operations.

```
Imported module: math
area_of_circle is a method definition
rd is a variable
  Its value is Constant(value=5)
ar is a variable
  Its value is Call(func=Name(id='area_of_circle', ctx=Load()), args=[Name(id='r', ctx=Load())], keywords=[])
```

Fig. 17. Generated explanation of the code in Fig. 16.

```
[{'Question': 'What is the name of the module that is imported?', 'context': 'Imported module:
math'}]
[{'Question': 'What is a method definition?', 'context': 'area is a method definition'}]
[{'Question': 'What is r?','context': 'r is a variable of type unknown'}]
[{'Question': 'What is Constant(value=5)?','context': 'Its value is Constant(value=5)'}]
Constant(value=5)'}]
[{'Question': 'What is a variable of type unknown?','context': 'a is a variable of type unknown'}]
[{'Question': 'What is the calculated area of the circle?','context': "'a' represents the calculated
area of the circle."}]
[{'Question': "What is the value of the call(func=Name(id='area', ctx=Load()),
args=[Name(id='r', ctx=Load())]?",'context': "Its value is Call(func=Name(id='area', ctx=Load()),
args=[Name(id='r', ctx=Load())], keywords=[])"}]
```

Fig. 18. Generated questions using the proposed model.

```
ZERO{}
ZERO{}
[{'Question': 'What is the area of the math.pi * radius?', 'context': 'area = math.pi * radius ** 2'}]
ZERO{}
ZERO{}
ZERO{}
```

Fig. 19. Generated question without describing the snippet.

## V. DISCUSSION

In this experiment, various code snippets were tested for translation using the proposed ontology and fed into the QuestionGen model to create open-ended questions. Table II outlines the test cases, the generated questions, and the difficulty level of the tested code. It is noticed that human evaluation of AQG results is more accurate than automatic assessments [10]. The validity of the generated code is rated on a scale of 1 to 5, where one represents the least validity and five indicates the highest validity. Difficulty is assessed based on script logic, with five denoting complexity and one representing simplicity. For instance, identifying variable assignments is relatively straightforward, while understanding inheritance is more challenging. Generating appropriate questions from sophisticated or advanced code snippets, such as those utilizing third-party libraries, still presents limitations. Composing accurate questions becomes increasingly tricky as code complexity and inter-line relationships grow. Consequently, further development is necessary to enhance outcomes. Addressing this need will lead to more advanced results. Nevertheless, this study introduces a new dimension to e-learning and supplements existing question-generation approaches that have proven effective in textual sources.

TABLE II. TYPES OF SYNTAX COVERED

| | Test case | Code level of difficulty | Generated question | Context | Generated question validity |
|---|---|---|---|---|---|
| a) | Variable declaration | 1 | What is the value of xfoo? | xfoo is a string variable and its value is 'foo' | 4 |
| b) | list declaration | 2 | 'What are the items in the list variable ab? | 'ab is a list variable and it has 2 items' | 5 |
| c) | Class declaration | 3 | What is a person? | Person is a class definition | 5 |
| d) | Instance and property initialization | 4 | What is a school an instance of? | 'school is an instance of the property' | 3 |
| e) | Variable initialization, instance initialization, property. | 5 | 'What is var1 an instance of?' | var1 is an instance of the Person class with name 'Jane' and age 25" | 4 |
| f) | Inheritance identification | 5 | Who does a student inherit from? | Student inherits from Person | 5 |
| g) | Libraries import | 4 | What is the name of the module that is imported? | Imported module: math | 4 |
| h) | Functions | 4 | What is a method definition? | area is a method definition | 3 |
| i) | Variable type | 4 | What is r? | 'r is a variable of type unknown' | 4 |
| j) | Functions result | 5 | 'What is the calculated area of the circle? | 'a' represents the calculated area of the circle. | 5 |

## VI. CONCLUSION

E-learning has become very popular recently, notably accelerated by the onset of the pandemic. One area that has gained considerable attention among researchers is the automatic generation of questions derived from learning materials. However, the predominant focus of existing efforts lies in generating questions from textual content. This work, however, concentrates on generating questions tailored for Python programming language learners derived explicitly from code snippets found in textbooks and course materials. Leveraging ontologies, this approach demands less computational resources, enhancing the scalability of the framework across diverse systems. The proposed framework harnesses ontological mapping, associating each syntactic element with its corresponding meaning and explanation. The process involves translating code into text and subsequently feeding this translated text into an AI-based model for question generation. It aims to alleviate the burden on educators and

reduce the repetition of the same questions for different groups of students. Moreover, the generated questions from code snippets serve to evaluate students' general understanding.

However, the proposed approach still has some limitations. The generation of questions relies solely on the QuestGen AI model, which can occasionally result in poorly phrased questions due to its AI nature. Additionally, the model might struggle to identify certain third-party libraries in complex code snippets. Hence, it represents an opportunity for future work to facilitate the insertion and categorization of concepts from all libraries. Finally, exploring alternative models such as GPT and expanding the framework to recursively process all imported libraries would enable a deeper understanding of complex syntactic structures. This enhancement would empower the ontology to explain code snippets better and generate more nuanced and fitting questions.

REFERENCES

[1] Y. Ham and B. Myers, "Supporting Guided Inquiry with Cooperative Learning in Computer Organization," in Proceedings of the 50th ACM Technical Symposium on Computer Science Education, Minneapolis, USA: ACM, Feb. 2019, pp. 273–279. doi: https://doi.org/10.1145/3287324.3287355.

[2] R. S. J. d Baker, A. T. Corbett, and V. Aleven, "More Accurate Student Modeling through Contextual Estimation of Slip and Guess Probabilities in Bayesian Knowledge Tracing," presented at the International Conference on Intelligent Tutoring Systems, in Lecture Notes in Computer Science, vol. 5091. Montreal, Canada: Springer Berlin Heidelberg, Jun. 2008, pp. 406–415. doi: https://doi.org/10.1007/978-3-540-69132-7_44.

[3] C.-Y. Chung and I.-H. Hsiao, "Investigating Patterns of Study Persistence on Self-Assessment Platform of Programming Problem-Solving," in Proceedings of the 51st ACM Technical Symposium on Computer Science Education, ACM, Feb. 2020, pp. 162–168. doi: https://doi.org/10.1145/3328778.3366827.

[4] C.-Y. Chung, C. Y. C. Edu, and I.-H. Hsiao, "From Detail to Context: Modeling Distributed Practice Intensity and Timing by Multiresolution Signal Analysis," presented at the 14th International Conference on Educational Data Mining, Virtual: International Educational Data Mining Society, Jul. 2021. [Online]. Available: https://educationaldatamining.org/edm2021/

[5] P. Brusilovsky, M. Yudelson, and I.-H. Hsiao, "Problem Solving Examples as First Class Objects in Educational Digital Libraries: Three Obstacles to Overcome Problem Solving Examples as Interactive Learning Objects for Educational Digital Libraries," J. Educ. Multimed. Hypermedia, vol. 18, no. 3, pp. 267–288, Jul. 2009.

[6] R. Cafolla, "Project MERLOT: Bringing Peer Review to Web-Based Educational Resources," J. Inf. Technol. Teach. Educ., vol. 14, no. 2, Apr. 2006.

[7] H. K. M. Al-Chalabi, "Evaluation of a Multi-Parameter E-learning System using Web 3.0 Technologies," presented at the 13th International Conference on Electronics, Computers and Artificial Intelligence (ECAI), Pitesti, Romania: IEEE, Jul. 2021, pp. 1–4. doi: https://doi.org/10.1109/ECAI52376.2021.9515191.

[8] H. K. M. Al-Chalabi and U. C. Apoki, "A Semantic Approach to Multi-parameter Personalisation of E-Learning Systems," presented at the International Conference on Modelling and Development of Intelligent Systems, in Communications in Computer and Information Science, vol. 1341. Sibiu, Romania: Springer International Publishing, 2021, pp. 381–393. doi: https://doi.org/10.1007/978-3-030-68527-0_24.

[9] P. Denny, A. Luxton-Reilly, and J. Hamer, "The PeerWise System of Student Contributed Assessment Questions," in Proceedings of the tenth conference on Australasian computing education, Wollongong, Australia, Jan. 2008, pp. 69–74. doi: https://dl.acm.org/doi/10.5555/1379249.1379255.

[10] N. Mulla and P. Gharpure, "Automatic Question Generation: A Review of Methodologies, Datasets, Evaluation Metrics, and Applications," Prog. Artif. Intell., vol. 12, no. 1, pp. 1–32, Jan. 2023, doi: https://doi.org/10.1007/s13748-023-00295-9.

[11] G. Kurdi, J. Leo, B. Parsia, U. Sattler, and S. Al-Emari, "A Systematic Review of Automatic Question Generation for Educational Purposes," Int. J. Artif. Intell. Educ., vol. 30, no. 1, pp. 121–204, Mar. 2020, doi: 10.1007/s40593-019-00186-y.

[12] J. Alshboul and E. Baksa-Varga, "A Review of Automatic Question Generation in Teaching Programming," Sci. Inf. Organ., vol. 13, no. 10, pp. 45–51, 2022, doi: 10.14569/IJACSA.2022.0131006.

[13] H. A. A. Ghanim, J. Alshboul, and L. Kovács, "Development of Ontology-based Domain Knowledge Model for IT Domain in E-Tutor Systems," Int. J. Adv. Comput. Sci. Appl., vol. 13, no. 5, pp. 28–34, 2022, doi: https://dx.doi.org/10.14569/IJACSA.2022.0130505.

[14] J. Alshboul, H. A. A. Ghanim, and E. Baksa-Varga, "Semantic Modeling for Learning Materials in E-Tutors Systems," J. Softw. Eng. Intell. Syst., vol. 6, no. 2, pp. 85–91, Aug. 2021.

[15] T. Urazova, "Building a System for Automated Question Generation and Evaluation to Assist Students Learning UML Database Design," University of British Columbia, 2022. [Online]. Available: https://open.library.ubc.ca/soa/cIRcle/collections/undergraduateresearch/52966/items/1.0413656

[16] S. Russell, "Automated Code Tracing Exercises for CS1," presented at the Computing Education Practice 2022, Durham, United Kingdom: ACM, Jan. 2022, pp. 13–16. doi: https://doi.org/10.1145/3498343.3498347.

[17] S. Sarsa, P. Denny, A. Hellas, and J. Leinonen, "Automatic Generation of Programming Exercises and Code Explanations Using Large Language Models," presented at the International Computing Education Research, Lugano, Switzerland: ACM, Aug. 2022, pp. 27–43. doi: https://doi.org/10.1145/3501385.3543957.

[18] M. Sh. Murtazina and T. V. Avdeenko, "The Constructing of Cognitive Functions Ontology," presented at the 14th International Symposium "Intelligent Systems, Moscow, Russia: Procedia Computer Science, 2021, pp. 595–602. doi: https://doi.org/10.1016/j.procs.2021.04.181.

[19] R. G. Golla, V. Tiwari, P. Chokhra, and H. Okada, "QuestGen AI." [Online]. Available: https://github.com/ramsrigouthamg/ Questgen.ai.

# An Enhanced Anti-Phishing Technique for Social Media Users: A Multilayer Q-Learning Approach

Asif Irshad Khan[1], Bhuvan Unhelkar[2]

Computer Science Department-Faculty of Computing and Information Technology,
King Abdulaziz University, Jeddah 21589, Saudi Arabi[1]
Muma School of Business, University of South Florida, Sarasota-Manatee Campus,
Sarasota, FL 33620, USA[2]

*Abstract*—As social media usage grows in popularity, so does the risk of encountering malicious Uniform Resource Locator (URLs). Determining the authenticity of a URL can be a highly challenging task, primarily due to the sophisticated attack structure employed by phishing attempts. Phishing exploits the vulnerabilities of computer users, making it difficult to discern between genuine and fraudulent URLs. To address this issue, a self-learning AI framework is required to warn social media users of potentially dangerous links. While several anti-phishing techniques exist, including blacklists, heuristics, and machine learning-based techniques, there is still a need for improvement in terms of detection accuracy. Hence, this study proposed a novel approach to combat phishing attacks using artificial neural networks, and the main aim is to create and validate the anti-phishing technique tool for detection accuracy. Initially, the URL data is collected, followed by preprocessing and then the analysis for malicious activity using the Logistic Bayesian Long Short-Term Memory model (LB-LSTM). The observed malicious URL features are extracted using multilayer Q-learning with the CaspNet and swarm optimization models. Analysis of these features enables the identification of a malicious URL, which is then removed, and the social media user is warned. The proposed technique attained a detection accuracy of 94.33%, Area under the ROC Curve (AUC) of 98.71%, Mean Squared Error (MSE) of 5.67%, Mean average precision of 88.67%, Recall of 98.67%, and F-1 score of 94.34%.

*Keywords*—*Multilayer Q-learning; anti-phishing model; social media users; machine learning; optimization; URLs; logistic Bayesian LSTM model*

## I. INTRODUCTION

Phishing is a deceitful online practice that employs social engineering and a particular strategy to deceive individuals using the internet to obtain their personal information or critical online data [1]. Businesses across various sizes and industries increasingly recognize the imperative of investing in anti-phishing solutions to protect their most valuable asset. Phishing is a fraudulent technique used to acquire sensitive customer information, such as classified data, via deceptive emails, counterfeit websites, dubious internet-based advertisements or promotions using forged Short Message Service (SMS) messages purporting to be from reputable service providers or online organizations, and other similar methods. Recent studies have shown that individuals with a sociable and trustworthy personality are more susceptible to

phishing schemes, mainly when actively participating in various social media platforms. A single instance of attack occurred on the social media platform Facebook, which enticed individuals to visit fraudulent websites designed to mimic the Facebook login page. The dissemination of the attack afterward extended to Facebook users via the promotion of acquaintances to access the hyperlink present on the original user's profile [2]. As an example, the security reports of Details and Patterns in 2017 revealed that a substantial sum of about $5 billion was stolen during the period spanning from October 2013 to December 2016. This financial loss was attributed to a W-2 phishing attempt, which affected a global population of over 24,000 individuals. W-2 phishing emails have recently been identified as one of the most dangerous kinds of phishing email scams, primarily due to their propensity to facilitate fraudulent tax filings and refund claims. Specialized deception encompasses using harmful code or crime ware, often installed on a personal computer or laptop, without the user's awareness [3]. Phishing may take several forms, including DNS poisoning, keyloggers, capturing meetings, damaging records, injecting information, etc. A new kind of malicious software called "ransomware" has emerged in recent years, allowing cybercriminals to run malicious code on a client's assets, locking them and demanding a payment "ransom" to unlock them. According to the CSO's announcement, 93% of all phishing emails nowadays are "ransomware." [4]. According to this study [5], the vast majority of victims of such crimes pay ransom demands quickly. Using Artificial Intelligence (AI) to glean information from customer files, postings, or social media and provide a timely warning is essential for fighting the phishing scourge. Though AI has promise, most current implementations need extensive client oversight, costly resources, and considerable management effort to be successful. In contrast to traditional detection and warning methods, current deep-learning systems are faster and more effective and don't need client mediation. The recent trend in Deep Learning (DL) based research has focused on extracting key highlights from text rather than conventional facts. The ability of deep-learning algorithms to effectively identify and predict concealed patterns within textual data is the primary reason why a traditional approach may struggle to detect or anticipate such patterns [6]. The highlights of this study are as follows:

- An innovative approach to mitigate phishing attacks targeting social media users through AI techniques.

- The development of a multilayer Q-learning model, in combination with CaspNet (Mul_Q_capsnet), and swarm optimization for the extraction of URL features.

- Utilizes a Logistic Bayesian LSTM model (LB-LSTM) to detect malicious behavior within social networking URLs.

- Provides warnings to social media users regarding potentially dangerous URLs.

Following is an outline of this paper. The related research is described in Section II, Section III outlines the materials and methods, Section IV describes performance analysis, Section V highlights the results and discussion of this research, and finally, final thoughts on the research work, like conclusion and future work, are listed in Section VI.

## II. RELATED WORKS

In today's environment, phishing poses the Internet with the most prominent challenge it must overcome. Many researchers have worked to create services that protect users from cyberattacks by detecting and blocking phished URLs using artificial intelligence (AI) and deep learning (DL) techniques [7]. Previous studies have developed and implemented two types of phishing detection systems: list-based and AI-based phishing identification systems. The list-based study identified many factors contributing to users' susceptibility to phishing attacks, including the lack of personal computer (PC) knowledge, inadequate understanding of security indicators, susceptibility to visual deception, and limited attention. Also, study in [8] investigated the persistence of successful phishing attacks despite ongoing efforts to mitigate associated risks. The findings of their investigations indicated that even when trained to identify phishing attacks, people were vulnerable at a rate of 53%. A study in [9] indicates that client PC data, client orientation, and the client's educational level are among the primary factors influencing whether clients open phishing emails.

A study in [10] introduced a system that generates a whitelist by logging the IP addresses of websites containing a login interface that a user visits. The system issues a warning if there is any inconsistency in the recorded website information when a user accesses a site. However, this approach raises concerns regarding websites users visit for the first time. In response to this challenge, work in Reference [11] has proposed a strategy to alert users on the web by maintaining an up-to-date whitelist of reputable sites. This strategy consists of two key components: a domain-IP address-matching module and extracting link attributes from the source code. This approach is further discussed in Reference [12]. A collaborative learning approach was employed for detecting phishing attacks in emails. Substituted selection methods were used to eliminate features unrelated to accuracy, achieving nearly 100% accuracy with just 11 selected features.

Study in [13] employed the Phi DMA approach, which used five layers: URL highlight layers, lexical layer, and whitelist layer, and accomplished an accuracy of 92%. In another review [14], the examination of phishing was identified through SVM. Author in [15] proposed an outrageous learning machine, a regulated AI calculation to determine spam accounts in SinaWeibo, Chinese miniature writing for a blogging site. Alberto et al. proposed an internet-based framework to filter comments posted on YouTube [16].

In study [17], the authors presented a structure for discovering dubious conversations on web-based gatherings using a coordinated help vector machine and particle swarm optimization methodology. The study by [18] focuses on detecting harmful URLs inside web-based social networks using client behavior analysis. The authors propose a research framework that investigates and detects social spam. This framework incorporates characteristics from URLs and online social networks (OSNs), emphasizing user profiles, postings, and URL attributes. The objective is to improve the accuracy of identifying harmful activities. A confirmation of the concept enhancer method was developed in [19], effectively used for the identification of bots, and in [20] identified spam in SMPs and involved the value of features in emphasizing a higher result collection of regulations. AI techniques require an environmental input to be adjusted and moved along.

The authors in [21] highlighted that the current apps, services, and systems are at risk of cyber-attacks like malware and software piracy because of the always-on nature of the Internet. These dangers threaten not just confidentiality but also safety. Malicious software like computer viruses, ransomware, scareware, and Trojan horses, as well as more traditional forms of software piracy like hard-disk loading, client-server overuse, and internet piracy, have the potential to wipe out critical data, resulting in reputational and economic damage. Reference [22] suggests companies can comprehensively enhance their operations by emphasizing roles, processes, individual actions, business strategies, business process modeling, quality assurance, cybersecurity, accountability, and big data.

In research [23], the author used a Neural Network (NN) to examine the blunder level of 4000 bogus and 4000 genuine pictures. With a solid achievement rate, a certified neural network has figured out how to group images as misleading or valid. Feature extraction extracted 17 features from 2500 phishing URLs from the PhishTank archive [24] and divided them into address bar-based highlights, unusual-based elements, and HTML and JavaScript-based highlights. Most parts were automatically separated from the URL and the page's source code without depending on third-party services. However, the WHOIS extracted the domain's age and DNS record [25]. The Alexa database retrieved the page's ranking [26]. Concurrently, the authors outlined an IF-ELSE rule and assigned a weight to each element. The weight of a feature was established by calculating the feature value concerning the total number of phishing links. Each segment's value could be either 1, 0, or 1, representing legitimate, suspicious, or phishing [27].

## III. MATERIALS AND METHODS

This section discusses novel techniques in the anti-phishing model using machine learning techniques for social media users and network optimization. The proposed architecture is shown in Fig. 1.

Fig. 1. Proposed Self-Learning framework for detecting malicious URLs.

## A. Data Pre-processing

The acquired data must undergo several preprocessing steps before entering each classifier to ensure that the algorithm interprets the information correctly and selects the best approach. One of these preprocessing activities involves organizing and cleaning the data. Formatting is crucial for presenting data in a format that classifiers can understand, such as converting data types into a text file or a tabular form. The cleaning process addresses missing values in a dataset, such as missing names or values in specific data fields. It involves setting properties determined manually or by the majority vote for matching values in different instances and even removing specific examples that could negatively impact the classifier's learning process. Additionally, cleaning details refer to the removal of personal information that could compromise the privacy of specific individuals.

## B. Tokenization

Breaking down text into its constituent parts, which can be words, sentences, or even individual characters, is known as tokenization, and the individual units are referred to as tokens. The objective here is to analyze text as a single unit. The list of tokens then becomes an interpretive response or serves as input for further sentence-based analysis. In languages where sentences are divided into segments and computer science, tokenization is crucial for reading text. Proper text tokenization is typically essential at the beginning of any text analysis process. All recognized text analysis methods rely on terms extracted from the dataset. To achieve this, a processor needs to tokenize the data. This can be straightforward when the text is in computer-readable formats. However, specific challenges may arise, such as handling punctuation marks. Other characters, such as parentheses, hyphens, and so on, also need to be managed.

## C. Logistic Bayesian LSTM Model (LB-LSTM)- based Malicious Activity Analysis

In logistic regression, the dependent variable typically takes on a binary form, which means it has only two possible values, like 0 or 1, true or false, or yes or no. This characteristic makes logistic regression well-suited for predicting the probability of an event belonging to one of two categories: success or failure. In this scenario, a sigmoid function is commonly employed to describe the connection between the predictor variables and the likelihood of the event happening. The sigmoid function yields

output values ranging from 0 to 1, effectively representing probabilities.

Consider a prototypical example with two predictors, A1 and B2. These predictors can be either constant values or binary variables, taking values of 0 or 1. The conversion likelihood W(A B) may be further divided into the acceptance probability A(A B) and the trial proposition probability T(A B), resulting in the equation W(A B) = T(A B) • A(A B). The likelihood of moving from state A to state B and the likelihood of moving in the other direction are related via Eq. (1):

$$P_A \cdot T(A \rightarrow B) \cdot A(A \rightarrow B) = P_B \cdot T(B \rightarrow A) \cdot A(B \rightarrow A) \quad (1)$$

For the likelihood for sample structure A to be equivalent to the Boltzmann weight by Eq. (2), the test plan and recognition likelihoods must be carefully selected.

$$P_A = \frac{e^{-\beta E_A}}{Z} \quad (2)$$

where, EA represents structure A's energy, since the conversion likelihoods are available using the proportion of probabilities, data for divider constant Z is unnecessary. Using Eq. (2), the detailed balancing requirement (3 may be rewritten as follows:

$$\frac{T(B \rightarrow A) \cdot A(B \rightarrow A)}{T(A \rightarrow B) \cdot A(A \rightarrow B)} = \frac{P_A}{P_B} = e^{-\beta(E_A - E_B)} \quad (3)$$

a process where all pairings of the states A, B satisfy the constraint T(A B) = T(B A). In the case of a design with N spins, for instance, it parallels picking one spin randomly from the matrix,: T(A B) = T(B A) = 1/N. Quickest if A(A B) or A(B A) is equivalent to 1, or if the larger of the two acceptance probabilities. 1 Padd is the chance of not adding an aligned spin. The complete balancing requirement may be expressed as Eq. (4).

$$\frac{T(A \rightarrow B) \cdot A(A \rightarrow B)}{T(B \rightarrow A) \cdot A(B \rightarrow A)} = \left(1 - P_{\text{add}}\right)^{m-n} \frac{A(A \rightarrow B)}{A(B \rightarrow A)} = e^{-\beta(E_B - E_A)} \quad (4)$$

Noticing that EA − EB = 2J (n − m), it follows that by Eq. (5):

$$\frac{A(A \rightarrow B)}{A(B \rightarrow A)} = \left[\left(1 - P_{\text{add}}\right)e^{2\beta J}\right]^{n-m} \quad (5)$$

Due to the constraint that a Bayesian network connected to Gi can have a maximum in-degree equal to n classes, there may be conditional probability tables with an exponential number of items in the nth category. However, this becomes impractical when dealing with issues that involve multiple types. Therefore, to address this challenge, it is necessary to reduce the maximum in-degree through the application of a structural learning technique. In real-world applications, we evaluate the Gi optimization as follows in order to shorten the arcs connecting classes to features arcs: $\mathcal{G}_i^* = \arg\max_{G,CGC} \log P(\mathcal{G} \mid \mathcal{D})$, where it's important to consider set inclusions among graphs in the arcs space, as specified by Eq. (6).

$$\log P(\mathcal{G} \mid \mathcal{D}) = \sum_{i=1}^n \psi_a[C_i, \text{Pa}(C_i)] + \sum_{j=1}^m \psi_a[F_j, \text{Pa}(F_j)] \quad (6)$$

Pa(Fj) represents Fj's paternities consistent with G, and Pa(Ci) defines Ci's parentages. Moreover, ψα is BDEu score

with equivalent model size α. For illustration, the score ψα[Fj ,Pa(Fj )] is $\sum_{i=1}^{|Pa(F_j)|}\left[\log\frac{\Gamma(\alpha_j)}{\Gamma(\alpha_j+N_{ji})}+\sum_{k=1}^{|F_j|}\log\frac{\Gamma(\alpha_{ji}+N_{jik})}{\Gamma(\alpha_{ji})}\right]$.

Finally, while j = P i ji, ji is equal to divided by the sum of the (joint) states of the parents of Fj and the number of states of Fj. All class variables share the same parents across all the graphs within the search space when considering a network associated with a specific class event, as the connections linking the class events are predetermined. This implies that the initial sum in Eq. (4) remains constant on the right side. Consequently, by concentrating solely on the attributes, we can achieve the optimization described in Eq. (3). A feature's parent set may be selected from any subset of C, which simplifies the problem into m distinct local optimizations. In reality, Eq. (7) by G establishes who Fj's parents are.

$$Bold\mathbf{C}_{F_j} = \arg_{Pa(F_j)\subseteq C}\max_a[F_j, Pa(F_j)] \qquad (7)$$

For each time j = 1, m. This is made practical by bipartite partition of class occasions as well as elements, yet much of the time, coordinated cycles might be found in the chart that expands every neighborhood score. Let k be the quantity of blend parts (i.e., the quantity of C qualities), with X being the arrangement of inquiry factors, Z being different factors. The marginal distribution of X may be calculated using Eq. (8) by adding together C and Z:

$$P(X=x) = \sum_{c=1}^k\sum_z P(C=c,X=x,Z=z)$$
$$= \sum_{i=1}^k\sum_z P(c)\prod_{i=1}^{\Pi}P(x_i\mid c)\prod_{j=1}^H P(z,\mid c)$$
$$= \sum_{c=1}^k\sum_z P(c)\prod_{i=1}^{|1}P(x_i\mid c)\prod_{j=1}^4\sum_{z_j}P(z_j\mid c) \qquad (8)$$
$$= \sum_{c=1}^k P(c)\prod_{i=1}^M P(x_i\mid c)$$

where past equality is valid since, for any j, j 1 zj P c z like this, it is easy to dismiss non-inquiry factors Z while figuring P(X = x), and calculation of P(X = x) takes O(|X|k), paying little heed to |Z|. Bayesian organization inference, conversely, is the most pessimistic scenario dramatic in |Z|. Restrictive probabilities, successfully assessed as proportions of peripheral probabilities, should likewise be considered. P(X=x,Y=y) = P(X=x,Y=y)/(Y=y ) The combination of trees, where every variable in each bunch is permitted to have one extra parent notwithstanding C, gives somewhat more extravagant model than naïve Bayes while as yet taking into consideration productive inference.

The computation involved is about the probability $(P_1^-, P_2 \dots P_n)$ when $X(x_1, x_2 \dots x_n)$ belongs to $C_1, C_2 \dots C_n$, with $P_j X(x_1, x_2 \dots x_n)$ being the probability when $X(x_1, x_2 \dots x_n)$ belongs to $C_j$. Then max $(P_1, P_2 \dots P_n)$ is demanded result.

$$P(C_j \mid x_1, x_2 \dots, x_n) = P(x1, x2, \dots, x_n \mid C)P(C_j) \qquad (9)$$

According to this formula, P(Cj) is the posterior probability that Cj includes text vector (x1,x2... xn) when the text to be categorized belongs to Cj and P(x1,x2... xn|Cj) is the prior probability when text belongs to Cj. As a result, max (P1,P2,...,Pn) is the highest value possible for the following Eq. (10):

$$\underset{x_t}{\arg\max}\, P(C_j \mid x_1, x_2 \dots, x_n) \qquad (10)$$

The qualities (x1,x2,...,xn) are assumed by Bayes to be independent of one another. The product of the probabilities for each attribute is then the joint probability. Thus, Eq. (11) is the final classification function.

$$\underset{x_t}{\arg\max}\, P(x1, x2, \dots, x_n \mid C)P(C_j) \qquad (11)$$

In this formula, $\frac{N(X_i=x_i,C=C_j)+1}{N(C=C_i)+M}$ is amount of text in training set that belongs to Cj, and N is total amount of text in training set. $P(x_i \mid C_j) = \frac{N(X_i=x_i,C=C_j)+1}{N(C=C_i)+M}N(X_i=x_iC=C)$ is the amount of text that has Cj's attribute xi, N(C=Cj) is amount of text that belongs to Cj, and M is text vector's dimension. Naive Bayes classifier's core design process, where two areas have been enhanced, is described above.

Typically, when presented with a microarray picture, we are unsure of which category it belongs to, hence the fairness principle demands that the text obtain the same prior probability for each group. Since there are differing amounts of text types in training sets, treating the prior probability differently is unjust and illogical. As a result, it makes sense to evaluate prior probability computation and apply the same prior probability instead. The classification function Eq. (12) shown below can therefore be obtained:

$$C_j \in C_i = 1\prod_{i=1}^n P(x_i \mid C_j) \qquad (12)$$

Omitting the computation of prior probability can significantly accelerate the calculation, but it does not affect the final grouping result because the maximum probability is required. The proposed model uses a two-valued constant to store the subsequent likelihood. Sometimes, the sentence that needs to be categorized is lengthy, which increases the sentence's dimension vector and decreases the subsequent likelihood. Additionally, reproducing the probabilities of all potentials may lead to inaccurate transmission. This can be corrected by reproducing the subsequent likelihood of each feature property. This will not affect the experiment's findings since what matters in the end is comparing probabilities between categories, and multiplying probabilities is logical. Initially, we expanded ten times, but throughout the research, we found that the subsequent likelihood would occasionally fall outside the range of the two-valued constant, significantly impacting the experiment's outcomes. We added an enlargement feature K to the function to reduce the impact of inaccurate propagation. The following optimization Eq. 13 and Eq. 14 are optimized to evaluate the load and bias values.

$$\min_{w,b,\xi}\frac{1}{2}ww^T + C\sum_{i=1}^n\xi_i \qquad (13)$$

$$\text{subject to} \begin{cases} y_i(w^T\phi(x_i)+b) \geq 1-\xi_i \\ \xi_i \geq 0, i = 1, \dots, n \end{cases} \qquad (14)$$

$\phi(x_i)$ – other dimension data points have been transferred. ξi- represent slack variable and these variables direct observations in the direction of the margin. Regularization is defined by C.

In this study, LSTM networks are utilized to create base models. Fig. 2 displays the construction of our base learner models using LSTM. Our deep LSTM network consists of multiple hidden layers, with a sequential input layer and two LSTM hidden layers added between the input and output layers. Each hidden layer consists of multiple memory cells and fully connected dense layers. We incorporate multiple LSTM and dense layers to construct a more precise, robust, expressive deep network for URL classification. The LSTM strategy takes URLs as input without requiring manual element extraction.



Fig. 2. Structure of the proposed LSTM model of the system.

The LSTM algorithm automatically processes the sequence of characters found in the URL. URLs are initially subjected to a tokenization process, where each unique character in the sequence is assigned a specific number, thereby converting the URL into tokens. These resulting token lists serve as inputs to the LSTM base models. The output layer makes predictions regarding the final outcome, specifically whether the URL is associated with phishing. We conduct experiments to determine the optimal number of LSTM units, the number of neurons in each dense layer, and the number of dense layers in each network. The models are fine-tuned by varying the number of epochs.



Fig. 3. Logical structure of the proposed LSTM model of the system.

Fig. 3 illustrates the proposed methodology. The first phase involves creating and training n LSTM models, each using a distinct subset of the training data. In the subsequent stage, a collection of records is supplied as input to each model for testing purposes. The models are tasked with predicting the

classification of each respective test record. The proposed ensemble LSTM model utilizes a voting approach for generating predictions. Essentially, the class assigned by the ensemble model is based on the majority of votes received from the individual models for a given test record.

For instance, if five LSTM models trained on distinct subsets of training data are provided with a particular test record, and four of them predict the record as "phishing" while the remaining one predicts it as "legitimate", the ensemble method would predict it as "phishing" using the voting approach since it was the majority prediction among the individual models. Algorithm of the proposed model is given below.

| **Algorithm 1:** Phishing Detection with LB-STM, CaspNet, and Swarm Optimization |
| --- |
| **1: Input:** Social media URL data |
| **2: Output:** Warning for potentially dangerous links |
| **3: Initialize:** |
| • URL data collection module. |
| • LB-LSTM model for malicious activity analysis. |
| • Multilayer Q-learning with CaspNet and swarm optimization models for feature extraction. |
| • Anti-phishing tool for social media users. |
| **4: Main Loop (URL Analysis):** |
| **5: for** all URL ∈ SocialMediaURLData **do** |
| **6:** Preprocess URL data: preprocessed_data ← Preprocess(URL) |
| **7:** Extract features using multilayer Q-learning, CaspNet, and swarm optimization: q_features ← MultilayerQlearning(preprocessed_data) caspNet_features ← CaspNet(preprocessed_data) swarmOpt_features ← SwarmOptimization(preprocessed_data) |
| **8:** Combine extracted features: Combined_features ← [q_features, caspNet_features, swarmOpt_features] |
| **9:** Analyze malicious activity using LB-LSTM model: maliciousness ← LB-LSTM(combined_features) |
| **10:** **if** maliciousness > Threshold **then** |
| **11:** Remove the URL and warn the social media user |
| **12:** **end if** |
| **13: end for** |

### D. Swarm Network Optimization

The Particle Swarm Optimization (PSO) is a system that utilizes particles traveling at a certain speed in a swarm or population to explore potential solutions and address each competitor arrangement. The system is made up of two stages: the preparation stage and the location stage. During the preparation stage, various legitimate and phishing websites are utilized to create and construct a sophisticated recognition model. Initially, 11055 websites are divided into twelve distinct categories based on the characteristics derived from the website's address bar. Additionally, six categories are established based on anomalies, five categories pertaining to HTML, and classifications dependent upon JavaScript.

Furthermore, the website's domain is used to classify seven distinct groups.

To determine the weightings to apply to specific website qualities, the suggested technique employs Particle Swarm Optimization (PSO). This practice is used due to the unique importance and varied contributions of each feature in identifying phishing websites. The PSO method accurately identifies phishing websites by computing the weighted sum of webpage properties. This practice guarantees that essential information is included in the artificial intelligence computations. In order to optimize the efficacy of artificial intelligence (AI) models, the technique of component biasing is used. This methodology assigns diminished weights to traits of lesser influence while attributing more significance to pivotal features. This stands in contrast to component selection techniques that completely neglect less relevant features.

It is anticipated that component weights will be represented as real numbers falling within the [0, 1] range. In accordance with the Particle Swarm Optimization (PSO) algorithm's recommended feature weighting, these values will indicate the relative importance of website characteristics. The number of features to which weights are assigned in the PSO is determined by the dimensionality of each particle.

Consider a scenario where 'n' feature weights are encoded in PSO particles according to the suggested PSO-based feature weighting method. In this setting, the 'nth' component and any supplementary features have a disproportionate weight, whereas the principal feature has less sway. However, the 'third' and 'n-1st' properties are deemed unnecessary when developing reliable forecasting models. A novel method for feature weighting, employing Particle Swarm Optimization (PSO), has been introduced. PSO entails a swarm of particles exploring diverse feature weight configurations. Each particle moves with random speeds and angles, characterized by a unique position representing the assigned weights for various features. These weights are encoded as real values within the range of zero to one.

Once the initial particle swarm is established, each particle's fitness is assessed based on the accuracy of feature characterization. PSO evaluates the well-being of each particle by initially determining the characterization accuracy and utilizing the weighted attributes of the particle for constructing the AI model from the training data. The objective of the PSO's health evaluation is to identify the best position for each element (pbest) as well as the best position for the entire swarm (gbest). If the values of interest surpass the individual health benefits of previous pbest and gbest, the current pbest and gbest will be updated. Each particle then appropriately adjusts its speed and position according to the updated pbest and gbest. This process is iterated until PSO optimizes the most prominent features. Ultimately, PSO yields the gbest, which contains some of the best feature weights achievable within the swarm.

Next, six commonly used AI algorithms are generated using a training dataset with features weighted through PSO's feature weighting process. The dataset, enriched with components weighted by PSO, is harnessed to construct various machine learning models, including Backpropagation Neural Networks, Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), Decision Trees (DT), Random Forests (RF), and Naïve Bayes classifiers (NB). These pre-trained models are developed and stored according to the specifications necessary for the effortless identification of new phishing websites.

A thoroughly evaluation is conducted using a recently collected dataset to gauge the effectiveness of the developed methods for detecting phishing websites. The primary objective is to determine the most vital components within the testing dataset. These elements include properties related to HTML and JavaScript, characteristics of the address bar, attributes based on geographical regions, and features linked to anomalies. The rigorous method of extracting features is crucial in developing accurate website models.

Subsequently, weights are applied to the characteristics in the testing dataset based on the optimum values generated from the PSO algorithm during the preparation step. This critical step significantly improves the accuracy of phishing website detection. PSO-weighted characteristics are added as input variables into the phishing site detection models created in the preliminary phase when considered necessary. These models are then used to determine if a website fits the requirements for being classified as a phishing site.

In conclusion, the performance of phishing site detection models incorporating proposed PSO-based feature weighting is evaluated and compared to standalone models.

## IV. PERFORMANCE ANALYSIS

The main goal of this research is to evaluate how effectively organizations can combat phishing attacks by utilizing artificial neural networks. Our objective is to design a machine-learning-based anti-phishing technique that can be easily integrated into social media platforms. To achieve this, we gather social media URLs and analyze them using the logistic Bayesian LSTM model (LB-LSTM) to identify any malicious activity. We use a combination of multilayer Q-Q learning with the CaspNet (Mul_Q_capsnet) and swarm optimization models to extract features that are indicative of malicious URLs. Any URLs that exhibit these features are then flagged and removed to ensure the safety of users.

Testing the proposed framework presents a significant challenge, as a globally recognized dataset is not readily available. To tackle this issue, we have undertaken the task of generating our own dataset. The dataset included in this study consists of 73,575 URLs, selected based on their considerable size and lack of device restrictions. The dataset contains 37,175 URLs classified as phishing and 36,400 URLs classified as legal. In the experimental configuration, we used a random assignment strategy to allocate 75% of the dataset for training purposes, while the remaining 25% was put aside exclusively for classification testing. To conduct a thorough assessment of the efficacy of our methodology, we used a range of measures, including accuracy, recall, and F-measure. It is important to note that each classification Algorithm was developed and tested independently for user-based features. In the same way, each classifier is trained and tested separately for content-based characteristics.

## A. Dataset

As previously discussed, in our efforts to obtain a reliable dataset for assessing the proposed framework, we faced challenges in finding one that fulfilled our standards. Thus, it was crucial to create a strong and extensive dataset. To accomplish this, we concentrated on gathering two distinct sets of URLs - legitimate and phishing - to ensure a well-rounded and thorough dataset. For the phishing URLs used in this research, our primary source was PhishTank (2018). However, it's worth noting that PhishTank does not provide a free dataset. As a result, we implemented a script to systematically acquire a substantial volume of rogue website addresses. At the same time, we collected authentic websites. The Yandex Search API proved to be invaluable in this regard (YandexXMLYandex Innovations, 2013).

To obtain a collection of web pages with minimal risk of phishing, we first created a specific query_word_list and then utilized the Yandex Search API to retrieve the top-ranked pages. This method was chosen due to the transient nature of malicious URLs, resulting in lower rankings by search engines. Our efforts resulted in a comprehensive dataset of 73,575 URLs, now available online for fellow researchers. This dataset includes 37,175 phishing URLs and 36,400 legitimate URLs, ensuring a well-rounded evaluation.

## B. Feature Analysis Based on NLP (Natural Language Processing)

By utilizing information preprocessing as detailed in previous sections, separating a few unmistakable features can be made simple. These features are extracted using Natural Language Processing (NLP) techniques that rely on the English language. For optimal efficiency, features are currently extracted based on the English language but can easily be adapted to any language. The selection and design of these features are minor issues, and most works focused on phishing detection use various feature lists based on their algorithms. Our chosen feature list primarily includes the need to parametrize the URL of the webpage, meaning that the text-based web address should be broken down into the words it contains. However, this process is not a simple task as a web address can contain several concatenated texts, making it difficult to find each word. To address this, the URL is parsed, and certain unique characters are considered, such as ("?", "/", ".", "=", and "and"), to rot the text.

The last character is particularly favored by attackers to convince victims that it is a legitimate webpage. Attackers use various tactics to deceive users, including using publicly-known brand names like Apple, Google, Gmail, Oracle, or specific keywords such as login, secure, account, server, etc., depending on the type of attack or targeted website. Therefore, in addition to the features proposed in previous literature, we defined several additional elements for detecting phishing websites. Although this number is not excessive, it is necessary to apply a feature reduction tool when using NLP, either alone or in combination with other techniques.

## C. Word Vectors

Converting words into vectors is a popular approach for identifying key features, such as text handling or text mining

techniques. Our system connects with the URL of a webpage, which is essentially a message composed of many words. Rather than manually modifying these words, a programmed vectorization technique is preferred. To accomplish this, we use the "StringtoWordVector" feature of Weka to convert each URL into a word vector. Once the linked vectors are obtained, the selected AI algorithm can easily utilize them. In the suggested framework, we tested 73,575 URLs and extracted 1,701 word highlights during the vectorization process.

To reduce the number of highlights, we utilized a component reduction system that employs the "CfsSubsetEval" method - an algorithm for element determination that utilizes the best first pursue technique. This lowering tool has reduced the number of necessary highlights from 1,701 to 102 in the rundown.

## V. RESULTS AND DISCUSSION

### A. Results

Finding a widely recognized dataset was one of the most difficult challenges we encountered when evaluating our suggested methodology. We couldn't find one, so we made our own. This dataset has 73,575 URLs and was selected owing to its massive size and absence of a test gadget restriction. The collection contains 37,175 phishing URLs and 36,400 trustworthy URLs. Our tests were conducted on a MacBook Pro with a 2.7 GHz Intel Core i5 CPU and 8 GB of 1867 MHz DDR3 RAM. We used Weka and numerous readymade libraries to test our system. We used 10-fold Cross-Validation and the default limit possible gains of all computations during the testing. In addition, each test set was run using seven distinct simulated intelligence algorithms. We created a confusion matrix for the tested learning algorithms and eventually found the optimal test type, whether NLP-based features, Word Vectors, or Hybrid. Table I shows the parametric study results of the suggested anti-phishing model in terms of Detection accuracy, AUC, MSE, Mean average precision, Recall, and F-1 score. Both the training and test data are analysed when the parameters are adjusted. Fig. 3 is a confusion matrix illustrating the identification of malicious URL detection.

Table I shows the parametric study results of the suggested anti-phishing model in terms of Detection accuracy, AUC, MSE, Mean average precision, Recall, and F-1 score. Both the training and test data are analysed when the parameters are adjusted. Fig. 4 is a confusion matrix illustrating the identification of malicious URL detection.

TABLE I. PROPOSED ANTI-PHISHING MODEL-BASED PARAMETRIC ANALYSIS

| Metrics | Training results | Testing results |
|---|---|---|
| Detection accuracy | 94.54 | 94.33 |
| AUC | 98.55 | 98.71 |
| MSE | 5.46 | 5.67 |
| Mean average precision | 89.09 | 88.67 |
| Recall | 98.54 | 98.67 |
| F-1 score | 94.54 | 94.34 |

(a) Training result-based confusion matrix.



(b) Test result-based confusion matrix.

Fig. 4. Proposed model-based malicious URL detection based on confusion matrix for (a) training results, (b) Test results.

The precision-recall and ROcurve result for the proposed model is shown in Fig. 5, as determined by both training and testing. These results are analyzed based on the true positive and false positive rates, which are determined by the model's confusion matrix.



(a) Training PR- curve.



(b) Testing PR- curve.



(c) Training ROC curve.



(d) Testing ROC curve.

Fig. 5. Training and testing result analysis based on PR and ROC curve.

The proposed analysis, utilizing the training and testing results is shown in Fig. 6.

Fig. 6. Proposed training and test result analysis.

From the above figure, the proposed technique attained Detection accuracy of 94.54%, AUC 98.55%, MSE of 5.46%, Mean average precision 89.09%, Recall of 98.54% and F-1 score of 94.54% for training results; for testing results proposed technique attained Detection accuracy of 94.33%, AUC 98.71%, MSE of 5.67%, Mean average precision 88.67%, Recall of 98.67% and F-1 score of 94.34%.

TABLE II.    ANALYSIS BETWEEN PROPOSED AND EXISTING TECHNIQUE

| Metrics | SVM_RNN | DBM_SAE | Proposed_ LB-LSTM_ Mul_Q_capsnet |
|---|---|---|---|
| Detection accuracy | 87.23 | 91.8 | 94.33 |
| AUC | 92.1 | 95.2 | 98.71 |
| MSE | 10.56 | 8.74 | 5.67 |
| Mean average precision | 81.45 | 86.23 | 88.67 |
| Recall | 91.34 | 95.22 | 98.67 |
| F-1 score | 87.66 | 92.33 | 94.34 |

Table II shows analysis of the anti-phishing model in terms of Detection accuracy, AUC, MSE, Mean average precision, Recall, and F-1 score.



Fig. 7. Comparative analysis for anti-phishing model.

Fig. 7 analysis is shown. The proposed technique attained a Detection accuracy of 94.33%, AUC of 98.71%, MSE of 5.67%, and Mean average precision of 88.67%, Recall of 98.67%, and F-1 score of 94.34%, while existing SVM_RNN attained Detection accuracy of 87.23%, AUC 92.1%, MSE of

10.56%, Mean average precision 81.45%, Recall of 91.34% and F-1 score of 87.66%; DBM_SAE attained Detection accuracy of 91.8%, AUC 95.2%, MSE of 8.74%, Mean average precision 86.23%, Recall of 95.22% and F-1 score of 92.33%.

## B. Discussion

The anti-phishing technique presented here has shown impressive results, highlighting significant progress in terms of detection accuracy, Area under the ROC Curve (AUC), Mean Squared Error (MSE), Mean Average Precision, Recall, and F-1 score. Upon closer analysis of the methodology and results, several important factors come to light. Although the proposed solution based on Natural Language Processing (NLP) is generally effective, there is a concern regarding its ability to detect pages with a single-space name, like "www.testbank.com." Phishing attacks frequently take advantage of differences in URLs, and the model's inability to address this particular situation could limit its efficacy in real-world scenarios. In a typical phishing attack, the page is designed to appear legitimate, and attackers attempt to conceal their extended URL by using unusual words to deceive customers. Customers who are already aware of phishing attacks tend to look for shorter URLs. Tests show that straight or probabilistic AI models like support vector machines and linear regression perform poorly overall. On the other hand, tree-based models significantly improve the identification of phishing URLs and produce highly effective and significant results.

The study recognises the ever-changing nature of phishing attacks, but it does not extensively explore how well the model can adapt to these evolving tactics. Phishing techniques are always evolving, so it's important to evaluate how well the suggested technique can adjust to new threats and consistently detect them accurately in the long run. The study does not thoroughly address the importance of user awareness and education in relation to the proposed anti-phishing tool. Although the model is effective, it is crucial to prioritize user education to combat phishing attacks effectively. Considering this aspect would offer a more comprehensive approach to cybersecurity. The success of the model is greatly influenced by the calibre and variety of the training data. If the training data lacks diversity or does not accurately reflect various phishing scenarios, the model may face challenges in effectively applying its knowledge to real-world situations.

Staying Ahead of Changing Threats: Phishing techniques are always changing, and attackers are constantly finding new ways to get around detection mechanisms. It is essential for the model to be able to adjust to new phishing trends and variations in order to ensure its long-term effectiveness.

Implementing deep learning models, especially those with intricate architectures such as Logistic Bayesian Long Short-Term Memory (LB-LSTM), can require significant computational resources. Deployment on platforms with limited resources, such as mobile devices, can present challenges. False Positives and User Experience: Excessively sensitive models can produce false positives, incorrectly identifying valid URLs as malicious. This may result in a subpar user experience and a decrease in trust in the system.

Striking a balance between achieving accurate detection and minimising false positives can be quite challenging.

## VI. CONCLUSION

Phishing attack is a serious threat to any organization, as it preys on individuals' independent decision-making. Responding promptly and effectively to phishing attacks is crucial for maintaining a secure business environment. Since employees are often the primary targets of phishers due to their unpredictable online behavior, sophisticated attackers know how to bypass logical reasoning and take a more deceptive approach. This paper investigates current strategies for detecting phishing web pages using machine learning. Furthermore, this study aims to enhance the efficiency and effectiveness of phishing datasets by employing feature selection methods. Feature selection is used to optimize an up-to-date phishing dataset and expedite the model-building process. The study utilizes the logistic Bayesian LSTM model and feature extraction to assess malicious behavior, with Multilayer Q-learning and the CaspNet (Mul_Q_capsnet) model with swarm optimization applied in the process. The proposed method achieves impressive results, including an F-1 score of 94.34%, an AUC of 98.71%, an MSE of 5.67%, an MPP of 88.67%, an RPP of 98.67%, and a detection accuracy of 94.33%. Future research could focus on changes in user behavior and evaluating the significance of account suspensions as a parameter by collecting new data once these regulations have reached a stable state. Additionally, researchers could categorize the various Arabic dialects used on social media platforms, and our methodology could be extended to other popular Online Social Networks (OSNs) such as Facebook and Instagram. We will further explore the ability of the proposed technique to handle large datasets and real-time scenarios. Discover various deployment strategies for seamlessly integrating the model across multiple social media platforms, ensuring comprehensive protection for all users.

## AUTHORS' CONTRIBUTIONS

Each author made valuable contributions to the conception, validation, formal analysis, and initial draft writing of the study. A.I.K. and B.U. were responsible for the methodology, investigation, and visualization, with B.U. also providing oversight and contributing to the writing, review, and editing. All authors have reviewed and approved the final version of the manuscript for publication.

## DATA AVAILABILITY STATEMENT

Data used for this article were collected by the research team and will be given to other researchers upon request.

## CONFLICTS OF INTEREST

The authors declare no conflict of interest.

## REFERENCES

[1]  J. D. Rosita P and W. S. Jacob, "Multi-Objective Genetic Algorithm and CNN-Based Deep Learning Architectural Scheme for effective spam detection," International Journal of Intelligent Networks, vol. 3, pp. 9–15, 2022.

[2]  B. Prabhu Kavin et al., "Machine Learning-Based Secure Data Acquisition for Fake Accounts Detection in Future Mobile Communication Networks," Wireless Communications and Mobile Computing, vol. 2022, pp. 1–10, Jan. 2022.

[3]  A. S. Alhassun and M. A. Rassam, "A Combined Text-Based and Metadata-Based Deep-Learning Framework for the Detection of Spam Accounts on the Social Media Platform Twitter," Processes, vol. 10, no. 3, p. 439, Feb. 2022.

[4]  R. Ghanem, H. Erbay, and K. Bakour, "Contents-Based Spam Detection on Social Networks Using RoBERTa Embedding and Stacked BLSTM," SN computer science, vol. 4, no. 4, May 2023.

[5]  I. H. Sarker, Y. B. Abushark, F. Alsolami, and A. I. Khan, "IntruDTree: A Machine Learning Based Cyber Security Intrusion Detection Model," Symmetry, vol. 12, no. 5, p. 754, May 2020.

[6]  R. Ghanem and H. Erbay, "Spam detection on social networks using deep contextualized word representation," Multimedia Tools and Applications, Jul. 2022.

[7]  A. Almomani et al.,"Phishing Website Detection with semantic features Based on Machine learning Classifiers-A Comparative Study," International Journal on Semantic Web and Information Systems, vol. 18, no. 1, Jan. 2022.

[8]  Z. Zhang, R. Hou and J. Yang, "Detection of Social Network Spam Based on Improved Extreme Learning Machine," in IEEE Access, vol. 8, pp. 112003-112014, 2020.

[9]  N. Ahmed, R. Amin, H. Aldabbas, D. Koundal, B. Alouffi, and T. Shah, "Machine Learning Techniques for Spam Detection in Email and IoT Platforms: Analysis and Research Challenges," Security and Communication Networks, vol. 2022, pp. 1–19, Feb. 2022.

[10] B. Prabhu Kavin et al., "Machine Learning-Based Secure Data Acquisition for Fake Accounts Detection in Future Mobile Communication Networks," Wireless Communications and Mobile Computing, vol. 2022, pp. 1–10, Jan. 2022.

[11] S. Rao, Anil Kumar Verma, and T. Bhatia, "Hybrid ensemble framework with self-attention mechanism for social spam detection on imbalanced data," vol. 217, pp. 119594–119594, May 2023.

[12] M. Alshehri, A. Abugabah, A. Algarni, and S. Almotairi, "Character-level word encoding deep learning model for combating cyber threats in phishing URL detection," Computers and Electrical Engineering, vol. 100, p. 107868, May 2022.

[13] Q. Zhang, Z. Guo, Y. Zhu, P. Vijayakumar, A. Castiglione, and B. B. Gupta, "A Deep Learning-based Fast Fake News Detection Model for Cyber-Physical Social Services," Pattern Recognition Letters, vol. 168, pp. 31–38, Apr. 2023.

[14] W. Khan and M. Haroon, "An unsupervised deep learning ensemble model for anomaly detection in static attributed social networks," International Journal of Cognitive Computing in Engineering, vol. 3, pp. 153–160, Jun. 2022.

[15] E. Dubasova, A. Berdashkevich, G. Kopanitsa, P. Kashlikov and O. Metsker, "Social Network Users Profiling Using Machine Learning for Information Security Tasks," 2022 32nd Conference of Open Innovations Association (FRUCT), Tampere, Finland, 2022, pp. 87-92.

[16] K. Hayawi, S. Saha, M. M. Masud, S. S. Mathew, and M. Kaosar, "Social media bot detection with deep learning methods: a systematic review," Neural Computing and Applications, Mar. 2023.

[17] M. Bhattacharya, S. Roy, S. Chattopadhyay, A. K. Das, and S. Shetty, "A comprehensive survey on online social networks security and privacy issues: Threats, machine learning-based solutions, and open challenges," Security And Privacy, Oct. 2022.

[18] M. Senthil. Raja and L. Arun. Raj, "Fake news detection on social networks using Machine learning techniques," Materials Today: Proceedings, Mar. 2022.

[19] V. Niranjani, Y. Agalya, K. Charunandhini, K. Gayathri and R. Gayathri, "Spam Detection for Social Media Networks Using Machine Learning," 2022 8th International Conference on Advanced Computing and Communication Systems (ICACCS), Coimbatore, India, 2022, pp. 2082-2088.

[20] P. Manasa et al., "Tweet Spam Detection Using Machine Learning and Swarm Optimization Techniques," in IEEE Transactions on Computational Social Systems.

[21] S. K. Sharma, B. Bhushan, B. Unhelkar, "Security and trust issues in internet of things: Blockchain to the rescue," CRC Press, 2020.

[22] B. Unhelkar, T. Gonsalves, "Artificial Intelligence for Business Optimization: Research and Applications," CRC Press, 2021.

[23] A. Mughaid et al., "A novel machine learning and face recognition technique for fake accounts detection system on cyber social networks," Multimedia Tools and Applications, vol. 82, no. 17, pp. 26353–26378, Jan. 2023.

[24] Y. A. Alsariera, V. E. Adeyemo, A. O. Balogun and A. K. Alazzawi, "AI Meta-Learners and Extra-Trees Algorithm for the Detection of Phishing Websites," in *IEEE Access*, vol. 8, pp. 142532-142542, 2020.

[25] C. Singh and Meenu, "Phishing Website Detection Based on Machine Learning: A Survey," *2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS)*, Coimbatore, India, 2020, pp. 398-404.

[26] M. A. El-Rashidy, "A Smart Model for Web Phishing Detection Based on New Proposed Feature Selection Technique," Menoufia Journal of Electronic Engineering Research, vol. 30, no. 1, pp. 97–104, Jan. 2021.

[27] B. B. Gupta, K. Yadav, I. Razzak, K. Psannis, A. Castiglione, and X. Chang, "A novel approach for phishing URLs detection using lexical based machine learning in a real-time environment," Computer Communications, vol. 175, pp. 47–57, Jul. 2021.

# ML-based Meta-Model Usability Evaluation of Mobile Medical Apps

Khalid Hamid[1], Muhammad Ibrar[2], Amir Mohammad Delshadi[3],
Mubbashar Hussain[4], Muhammad Waseem Iqbal[5], Abdul Hameed[6], Misbah Noor[7]

Department of Computer Science, Superior University, Lahore, 54000, Pakistan[1, 6]
Department of Computer and Mathematical Sciences, New Mexico Highlands University, Las Vegas, USA[2, 3]
Department of Computer Science, University of Gujrat, 54000, Pakistan[4]
Department of Software Engineering, Superior University, Lahore, 54000, Pakistan[5]
Department of Languages & Communication, UniSZA, Gong Badak Campus, Universiti Sultan Zainal Abidin, Malaysia[7]

*Abstract*—Mobile medical applications (MMAPPs) are one of the recent trends in mobile trading applications (Apps). MMAPPs permit users to resolve health issues easily and effectively in their place. However, the primary issue is effective usability for users in maps. Barely any examination breaks down usability issues subject to the user's age, orientation, trading accessories, or experience. The motivation behind this study is to decide the level of usability issues, concerning traits and experience of versatile clinical clients. The review utilizes a quantitative technique and performs client try and hypothetical insight through the survey by 677 members with six distinct assignments on the application's point of interaction. The post-try review is finished with concerning members. The Response surface method (RSM) is used for perceptional and experimental designs. In each case, participants are divided into 13 runs or groups. Experimental groups are involved after checking the perceptions about theoretical usability for different attributes according to the usability model through the questionnaire. The difference is recorded between the perception of users about usability (theoretical usability) and actual performance for usability. The study analyzed through Analysis of variance (ANOVA) that there is a need to improve mobile medical applications but it is also recommended to minimize the gap between the perception level of laymen and the actual performance of IT literate users in context with usability. The experimentation measures the tasks usability of various mobile medical applications concerning their effectiveness, efficiency, completeness, learnability, memorability, easiness, complexity, number of errors and satisfaction. Every design model also produces a mathematical expression to calculate usability with its attributes. The results of this study will help to improve the usability of MMAPPs for users in their convenient context.

*Keywords*—*ANOVA; completeness; efficiency; effectiveness; perceptional usability; response surface methodology; actual usability*

## I. INTRODUCTION

A software program called a mobile application is created to operate on a mobile device, like a smartphone or tablet. The use of smartphones has expanded in many environments like banking, education and gaming including healthcare with many potential and real-life benefits. The younger generation becomes technically competent medical professionals. According to statistics, more than 36% of the world's population was using smartphones in 2018, up from about 10% in 2011. As the 2020 analysis shows, one of the Asian countries has approximately 82 million Internet customers, and the industry will exceed US$10 billion. In third-world countries, the proportion of individuals who own a mobile phone is even higher (75%) [1].

In January 2021, there were 61.34 million online customers in third-world countries. As of January 2021, the country has 173.2 million mobile associations. From January 2020 to 2021, the quantity of cell phone users in the nation expanded by 6.9 million (+4.2%). January 2021. In January 2021, the quantity of public compact affiliations is identical to 77.7% of the complete populace [2]. From a usability point of view, medical institutions have introduced the use of Internet technology to meet the needs of patients and improve their services, but this phenomenon is still in its infancy [3]. The most popular Internet service is the mobile medical application. Usability plays a significant part in the product improvement process. In the field of Human-Computer Interaction (HCI), the most generally acknowledged meaning of usability is that proposed in ISO 9241-11: "the degree to which a client can utilize an item to accomplish exact goals with efficiency, effectiveness and satisfaction in a setting of determining use ". Then again, in the field of software development, the most generally acknowledged meaning of usability is that proposed in ISO 9126-1 [4] "the capacity of the product item to be perceived, learned worked and alluring to the user when utilized under determined conditions ". Starting here in view, usability is considered a particular element that influences the nature of a product item; it doesn't be guaranteed to infer client interaction with the framework since it is very well may be estimated as "consistency to the determination". After indulging the efficiency, effectiveness, learning ability, reliability, safety, error-freeness, enjoyment and other factors of MMAPPs, it is of great benefit to provide a new and convenient way for the online mobile medical field [5]. Most of the users feel uncomfortable and dissatisfied with online medical treatment and consultation because MMAPPs are may not user-friendly. This is a new concept in developing countries such as Asian countries. Usability is also described as "the degree to which a specified customer can use a product to achieve specified goals in a pre-defined usage setting with feasibility, productivity, and realization" (ISO 9241-11, 1998) [6].

Health is wealth; this is an important lesson, especially considering the recent coronavirus outbreak. This deadly virus spread so fast that the entire city was abandoned to maintain it. This virus tells us that the study needs to connect the health system, and need to be more technically focused on identifying and treating such diseases without going outside in public places with the help of MMAPPs. Most of the functions performed by mobile medical apps are checkup waiting time, online patient evaluation, feedback, medical history, self-medication, first aid information, guidance to reach the hospital, simple payment method/simple payment, run-time diagnosis after entering symptoms, and variety of input methods (as this is the case of the patient), check the availability of specialist through cloud computing, sample collection facility, find a doctor, and set an appointment for a doctor. The most irritating thing is the hanging tight for specialists or the clinical benefits. According to the occupied cycle, the people groups can't bear the cost of hanging tight for the specialist or clinical benefits. Giving this information will show on the profile of the specialist how long you should hold back to get inspected or get any clinical benefit according to my perspective this one is the more charming element. MMAPs can achieve this after providing all functionalities [7].

This paper discusses the usability of mobile medical apps, reviews the literature, and discusses usability models, features, selections, and user-centered models. It includes research methodology, questionnaire selection and sampling, Central Composite Design I and RSM Experimental Design II results, compares the results, analyzes them and concludes with implications.

*A. Contributions*

- The study provides the ANOVA-based usability evaluation of mobile medical apps.

- The study calculates the perceptional usability and actual usability of mobile medical apps.

- The study provides the usability calculation formula with nine attributes of usability which can be used to calculate the usability of any of the mobile medical apps.

## II. LITERATURE REVIEW

The study explains the speedy growth of mobile users, there's an outstanding boom in mobile software users. Therefore, the preservation of cellular users and producers of mobiles increases processing power, storage, functionalities and offerings. Now there may be a task for builders, software program engineers and interface designers to play their element in context with usability to retain cell utility users. The motive is that every category of cellular packages like enterprise apps, schooling apps, leisure apps, clinical apps, Travel apps, software apps and social media apps has its own practical and non-functional needs [8].

The usability analysis of mobile apps is executed on the idea of the four maximum popular attributes which can be efficiency, effectiveness, usefulness, and accuracy. Analysis of this study explains that usability evaluation of some other

cellular utility may be carried out with the help of these four characteristics only [9].

The given examination evaluation suggests that social elements are extra effective in the reputation and usability of cellular packages. The have a look at is restrained due to the fact carried out only extracted facts of crowd sourced web packages. Due to human computations and tiny date units from only new jobs, validity and verification are sometimes jeopardized [10].

Its first degree used a user-targeted layout for customer's duties; the second degree examined usability with laboratory settings and third level usability assessment was performed in a real-world environment. As a result, it offers many usability evaluation strategies. Mobile utility builders can select first-class one or extra in keeping with the scenario [11].

The research evaluated the usability of mobile apps using 36 criteria to create applications that are centred around people. This look used three methods for reliable assessment; the methods are QUIM, mGQM, and GQM. According to participating specialists, the effects performed from this assessment are dependable and validated with the help of metrics [12].

The paper explains that cell applications need extra attention as compared to large display computers like laptops or desktop systems since mobiles have a small screen with ongoing warnings. It would be ideal for it to be simple and more clients lovely. A limit of the cell programs is assessed in a usability setting with ascribed productivity of one hundred%, adequacy of 96%, and pride of 87%, yet memorability, learnability, straightforwardness and mental burden are not assessed [13].

As per an examination, clients of portable and its programs developing quickly as there are 4.57 billion cell phone users. Those buyers use 175 billion projects/. This investigation presents the UCD form with usability credits as viability, execution, pride, understandability, blunders, and availability [14].

According to the report, most developers don't pay close attention to usability factors like accessibility and learnability most of the time. Mostly smart software engineers and interface designers are involved with effectiveness and satisfaction however they are not with user-side error safety Hamid et al. [15]. It also discusses 27 problems of usability, suggestions for them, and guidelines for those issues so that it will be beneficial for developers and researchers [16].

This study investigates the usability evaluation process of mHealth apps using a Systematic Literature Review. Results show that a mixed-method approach can improve reliability and satisfaction. The study encourages developers to design more user-friendly applications, especially for older adults and novice users, to improve the effectiveness of mHealth apps [17].

The study presents a methodological approach for developing a usable mHealth application using a three-level stratified health information technology usability evaluation framework. The methodology includes a card sorting technique

for user-task guidance, end-user testing and heuristic evaluation with experts, and real-world evaluation after a three-month trial. The case study illustrates the use of these methodologies. The three-level usability evaluation was used to explore user interactions, refine app content, and use a stratified health IT usability evaluation framework for mHealth app design, development, and evaluation, providing methodological recommendations for future studies [18].

## III. METHODOLOGY

The motivation behind this exploration is to investigate and foster comprehension of basic worries blunder-free, disappointment-free, more usable and best executable m-medical applications because the patient can't afford any misconception due to horridness [19]. This Idea will be founded on quantitative request. It is the orderly experimental examination of detectable peculiarities employing measurable, numerical, or computational techniques [20]. This examination contrasts client thinking with comprehending the convenience of versatile medical applications with specialists thinking based on gathered information and proposes the best technique and ideas for improvement of MMAPPs later investigated over chronicled data, and standards of conduct and made accessible to the exploration community [21].

### A. Usability Model

Model Produced with the help of the following attributes under consideration.

### B. Awareness (Interestingness)

The rising fame of mHealth is a promising and open door for torment self-administration. Versatile applications can be effortlessly grown, however understanding the plan and usability will result in applications that can hold more clients. This exploration targets recognizing, breaking down, and orchestrating the present status of the specialty of (a) the planned approach and (b) usability appraisal of agony the executive's portable applications [22].

### C. Complexity

Complexity analysis based on screenshots of the user interface in addition to interaction information, textual content size, font, language, or character set, homogeneous background and contrasting color without the want to get entry to the source code of the utility. In contrast, applications provide easy navigation and without blind flow [23].

### D. Easiness

An important concept that illustrates how well users can use a mobile application is the ease of use. Design engineers define specific KPIs for each project like "Clients should be able to tap Find within three seconds of reaching the point of interaction on the application interface" and "usability should be streamlined while providing the greatest usefulness and considering business constraints." [24].

### E. User Satisfaction

Fulfillment can be accomplished in three ways. As a matter of first importance, interface text or content should esteem the patient in setting with the significance of the patient and accomplishing the objective of the patient through the application. Furthermore, the point of interaction should direct the patient through the task for which the individual in question utilizing it [25]. Thirdly various assignments for finishing an exchange/accomplishing an objective ought to be well organized [26].

### F. Efficiency

The productivity of portable medical applications is assessed three correspondingly, as a matter of first importance either concerning application plays out the particular undertaking totally, precisely and brief time frame. Also, either concerning application load and login or logout in the brief time frame as indicated. Thirdly, whether the unsettling application is viable to different mobiles and human-PC association aptitudes [27].

$$\text{Time} - \text{Based Efficiency} = \frac{1}{N} \sum_{i=0}^{N} \frac{n_{ij}}{t_{ij}} \quad (1)$$

where, N= No. of Jobs, R= No. of Contestants, $n_{ij}$ = Job I's resulted by j's participant, and $t_{ij}$ = Participant I's time to Complete a j Task.

Eq. (1) calculates the time-based efficiency with the help of each total number of tasks completed by each participant in a specific time and divided by the number of jobs.

### G. Effectiveness

Before you start to arrange Effectiveness is estimated with the assistance of consistent appearance applications, either interface configuration has significant choices and fastens more noticeable, lucid and simple to get to. Furthermore, either client moves around various choices effectively and sensibly to explore versatile medical applications. Consequently, the viability is a mix of Logical appearance and navigation of the UI of versatile apps [28].

$$\text{Effectiveness} = \frac{Total\ Number\ of\ Tasks\ Completed\ Successfully}{Total\ Number\ of\ Tasks\ Undertaken} \times 100 \quad (2)$$

Eq. (2) calculates the effectiveness of different medical using programs that divide the total number of activities completed by the total number of tasks attempted and multiply the result by 100.

### H. Memorability

The idea of memorability, from the usability point of view, is that a client can leave a program and when the person gets back to it, recall how to get things done in it. Memorability is significant generally because clients may not be utilizing your application constantly. It is easy to recall a task which is previously performed and reconnect the user after a long time [29].

### I. Learnability

Learnability property signifies "How simple is it for the people to figure out how to utilize the framework". It tends to be accomplished assuming our product point of interaction is basic and has routine likenesses to the next application. People are not working any harder than needed to utilize innovation and try to avoid absolutely special software as individuals gained from past experiences. Various people have distinctive

trouble levels; it is likewise an observable highlight to accomplish/assess the learnability normal for portable medical apps [30].

### J. Completeness

Completeness means checking the application interface for style, buttons, navigation and task completeness, etc. As a result of the backward point of view of the UI versus the prerequisites, presented, it can re-decipher task culmination also. The inquiry assumes there are relations in the model, i.e., i.e., rules which oversee changes between states. In legitimate dialect, it needs to have the option to inquire as to whether the framework is finished. The difficulties of these rules closely match the ones which relate to task plumpness [31].

### K. Selection of Medical App Features

The study initiates a systematic study of the characteristics of more usable mobile medical apps and their impacts on the cyber world, and medical industry. The study compares different mobile medical apps for seven features 1) Find a Doctor, 2) Set an Appointment, 3) Sample Collection Facility, 4) Medical History, 5) Feedback and 6) Online Patient Evaluation, and proposes the best mobile medical apps with the use of a questionnaire instrument, the post-test is also carried out. Verification and validation of results have been carried out based on real-time data [32].

### L. User Center Approach

The User Centre Model (UCD) is a research technique that improves versatile applications by reinforcing their convenience and decreasing expense as seen in Fig. 1. The fundamental objectives of the UCD model are fulfillment, essential, learnable, compelling, proficient, and adjustable design or interface for the users. In this model collect the requirements from users, then develop designs accordingly through RSM, calculate usability attributes of perceptional usability, and IT User's usability and combine usability with the help of experiments on 13 runs or groups of users. The study evaluates and compares the results of perceptional usability, IT user usability and combined usability. After the analysis study produced coded equations or formulas for calculating usability.



Fig. 1. User center model.

### M. Reason for the Model Used

The study works with nine attributes of the usability for evaluation of mobile medical apps. The study provides the perceptional-based evaluation of usability version and actual means performance-based evaluation of usability and both are evaluated through ANOVA. At the end compare results of both evaluations.

## IV. SELECTION AND SAMPLING

The usability testing was driven at various university campuses and after filtration 677 individuals enlisted for the examination performed for twelve highlighted MMAPPs regarding interface for 6 features due to the availability of relevant participants and apps. The study divided the participants into 13 runs or groups starting from run 1 up to run 13. Each group was assigned several tasks according to RSM based design model for assessing and evaluating the usability of each feature concerning effectiveness, efficiency, learnability, memorability, completeness, easiness, complexity, number of errors and satisfaction [33].

### A. Questionnaire

The questionnaire is developed with the help of the System Usability Scale (SUS) and Post-Study System Usability Questionnaire (PSSUQ) and gathers data about nine attributes of usability from 677 participants. In this study five points grading scale is used from strongly agree to strongly disagree on the other hand mid-point is agree. The study assessed the ease of use, ease of learning, simplicity, effectiveness, efficiency, ease of memorable, awareness, completeness, information and the user interface [34].

There are two design models are applied in the study for calculating perceptional usability and combined usability, after applying the questionnaire and performing the experiment respectively.

## V. RESULTS AND DISCUSSION

The RSM (Response Surface Methodology) technique is used in this study to validate the usability model.

### A. Central Composite Design I (Perceptional Usability)

Table I represents different attributes of usability, their effects on usability and their response to perceptional usability

Table II shows variables, and their levels like minimum values, maximum values and mid values for all attributes of usability defined in a usability design.

Fig. 2 shows that validation of the model is done with the help of relationships and the effect of different attributes on usability given in the diagram. As seen from the above figure almost all the attributes of the concerning model affect usability. Some attributes have a greater effect and few have a little effect which is also shown in Eq. (3).

Table III shows the ANOVA model regression coefficient and analysis of variance which is significant and the lack of fit insignificant as required for the validation model.

TABLE I. ANOVA-BASED FACTORS AND THE RESPONSE OF PERCEPTIONAL USABILITY

| Run | Factor 1 A: Awareness | Factor 2 B: Complexity | Factor 3 C: Easiness | Factor 4 D: Satisfaction | Factor 5 E: Efficiency | Factor 6 F: Effectiveness | Factor 7 G: Memorability | Factor 8 H: Learnability | Factor 9 J: Completeness | Response 1 Perceptional Usability |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2.1 | 2 | 2.2 | 2 | 2.6 | 2.2 | 2 | 2.4 | 2.3 | 55 |
| 2 | 2.33 | 1.87 | 2.53 | 2.67 | 2.6 | 2.47 | 2.2 | 2.6 | 2.27 | 59.81 |
| 3 | 2.26 | 1.87 | 2.55 | 2.45 | 2.41 | 2.53 | 2.24 | 2 | 2.09 | 58.07 |
| 4 | 2.42 | 1.95 | 2.52 | 2.58 | 2.52 | 2.57 | 2.29 | 2.55 | 2.02 | 59.01 |
| 5 | 2.43 | 1.9 | 2.5 | 2.57 | 2.73 | 2.37 | 2.67 | 1.83 | 2.57 | 59.91 |
| 6 | 2.23 | 2.1 | 2.46 | 2.43 | 2.57 | 2.34 | 2.2 | 2.6 | 2.11 | 57.86 |
| 7 | 2.24 | 2.1 | 2.53 | 2.42 | 2.39 | 2.51 | 2.2 | 2.48 | 2.07 | 58.01 |
| 8 | 1 | 2 | 2 | 2 | 3 | 1 | 2 | 4 | 0 | 47.22 |
| 9 | 3 | 2 | 4 | 4 | 4 | 3 | 3 | 4 | 3 | 83.33 |
| 10 | 2.52 | 2.16 | 2.8 | 2.88 | 2.72 | 3.04 | 2.56 | 2.72 | 1.92 | 64.78 |
| 11 | 3.08 | 2.08 | 2.84 | 3.12 | 2.8 | 2.84 | 2.48 | 2.84 | 1.72 | 66.11 |
| 12 | 2.42 | 1.95 | 2.52 | 2.58 | 2.52 | 2.57 | 2.29 | 2.55 | 2.02 | 59.01 |
| 13 | 3.08 | 2.08 | 2.84 | 3.12 | 2.8 | 2.84 | 2.48 | 2.84 | 1.72 | 66.11 |

TABLE II. LEVELS OF INDEPENDENT VARIABLES

| Symbol | Independent Variables | Minimum Value | Mid Value | Maximum Value |
|---|---|---|---|---|
| A | Interesting (Awareness) | 1 | 2.38 | 3.08 |
| B | Complexity | 1.84 | 1.97 | 2.16 |
| C | Easiness | 2 | 2.64 | 4 |
| D | Satisfaction | 2 | 2.68 | 4 |
| E | Efficiency | 2.39 | 2.75 | 4 |
| F | Effectiveness | 1 | 2.47 | 3.04 |
| G | Memorability | 2 | 2.36 | 3 |
| H | Learnability | 1.83 | 2.73 | 4 |
| J | Completeness | 0 | 1.99 | 3 |



Fig. 2. Relationships and effects of attributes on perceptional usability.

TABLE III.  RESULTS OF REGRESSION COEFFICIENTS AND ANALYSIS OF VARIANCE FOR ONE RESPONSE VARIABLE BY ANOVA

| Source | Intercept ($\beta_o$) | A-Awareness | B-Complexitys | C-Easiness | D-Satisfaction | E-Efficiency | F-Effectiveness | G-Memorability | H-Learnability | J-Completeness | *P*-Value | F-Value | $R^2$ | Adj. $R^2$ | Lack of Fit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Theoretical-Usability | 59.36(Significant) | 4.91 | -1.13 | 8.12 | -1.69 | 11.09 | 7.92 | -0.72 | 1.42 | 1.75 | <0.0001Significant | 263.13 | 0.978918 | 0.975198 | 0.8551(Insignifican) |

## B. Analysis via Perceptional Design I

As seen in Fig. 3. Under half of the mobile patients have high and extremely high fulfillment levels on versatile medical applications. Over 59% of people said that mobile medical apps are too complex. Complexity is inversely related to usability. The vast majority of IT users and doctors said that current MMAPPs are successful, productive and easy to use but a small number of general users feel much easier. The graph shows that easiness is directly related to usability. The majority of the MMAPP's users think that these applications are satisfactory which is related to usability as seen from the graph. As seen in Fig. 1. little more than 50% of mobile medical application users think that these apps are efficient. The remaining of them think that more efficiency is needed. Efficiency is directly linked with usability as seen from the graph. 50% think effectiveness is necessary but the remaining do not think so. The graph shows effectiveness is a little bit related to usability, learnability and memorability are also directly linked with usability but have little effect on it.

$$Perceptional\ Usability = 59.36 + 4.19A - 1.13B + 8.12C - 1.69D + 11.09E + 7.92F - 0.72G + 1.42H + 1.72J \quad (3)$$

Eq. (3) represents the perceptional usability deduced from the ANOVA model after analysis. This is a general proposed formula of perceptional usability that can be used for the calculation of any sample of the study.

## C. Experimental Design

This test was overseen on university campuses where all members were grown-ups. The individuals were facilitated to perform 13 activities, for instance, the tasks were organized and executed for class length in the college. The ordinary task completion time was eight minutes [35].

## D. RSM Experimental Design II (Experimental Usability)

Table IV represents different attributes of usability, their effects on usability and their response to combined usability. There are 13 runs in which the specific number of participants is included according to the RSM design model from which particular tasks are performed to calculate and evaluate the attributes like effectiveness, efficiency and satisfaction. On behalf of these attributes, usability is calculated and the equation of combined usability is deduced.

Table V shows variables that represent the attributes of combined usability and their levels like minimum values, maximum values and mid values. This table also represents the standard deviation faced by given attributes in an RSM model.

TABLE IV.  ANOVA-BASED FACTORS AND RESPONSE AS ACTUAL USABILITY

| Run | Factor1 Effectiveness | Factor2 Efficiency | Factor3 Satisfaction | Response1 Usability |
|---|---|---|---|---|
| 1 | 82.93 | 75 | 50 | 69.31 |
| 2 | 57.63 | 83.33 | 66.67 | 69.21 |
| 3 | 85.44 | 84.12 | 61.18 | 76.91 |
| 4 | 81.71 | 84.12 | 63.55 | 76.46 |
| 5 | 58.46 | 63.33 | 63.33 | 61.71 |
| 6 | 72.33 | 70.14 | 60.71 | 67.73 |
| 7 | 80.03 | 81.74 | 61.18 | 74.32 |
| 8 | 75 | 75 | 50 | 66.67 |
| 9 | 50 | 100 | 100 | 83.33 |
| 10 | 69.32 | 72 | 72 | 71.11 |
| 11 | 65.9 | 72 | 71 | 69.63 |
| 12 | 81.59 | 87.06 | 63.55 | 77.4 |
| 13 | 59.47 | 60 | 70 | 63.16 |

TABLE V.  LEVELS OF INDEPENDENT VARIABLES FOR EXPERIMENTAL DESIGN

| Symbol | Independent Variables | Minimum Value | Mid Value | Maximum Value | Std. Dev. |
|---|---|---|---|---|---|
| A | Efficiency | 50 | 74.99 | 85.44 | 11.15 |
| B | Effectiveness | 60 | 87.3 | 100 | 10.27 |
| C | Satisfaction | 50 | 65.63 | 100 | 11.89 |
| R1 | Usability (Actual) | 61.71 | 75.97 | 83.33 | 6.14 |



Fig. 3.   Relationship and effects of attributes on combined usability.

TABLE VI. RESULTS OF REGRESSION COEFFICIENTS AND ANALYSIS OF VARIANCE FOR THREE RESPONSE VARIABLES

| Source | Intercept ($\beta_o$) | A-Efficiency | B-Effectiveness | C-Satisfaction | AB | AC | BC | $A^2$ | $B^2$ | $C^2$ | $P$-Value | F-Value | $R^2$ | Adj. $R^2$ | Lack of Fit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Actual-Efficiency | 7.12(Significant) | 1.13 | 1.12 | 1.12 | -0.13 | -0.11 | -0.12 | -0.059 | -0.045 | -0.053 | <0.0001 | 1.014E+008 | 1.0000 | 1.0000 | 0(Insignifican) |

Fig. 3 shows that validation of the model is done with the help of relationships and the effect of different attributes on combined usability given in the diagram. As seen from the above figure almost all the attributes of the concerning model affect usability. Some attributes have a greater effect and few have a little effect which is also shown in Eq. (4).

Table VI shows the ANOVA model regression coefficient and analysis of variance which is significant and the lack of fit insignificant as required for the validation model.

$$\text{Combined Usability} = 7.12 + 1.13A + 1.12B + 1.12C - 0.13AB - 0.11AC - 0.12BC - 0.059A^2 - 0.045B^2 - 0.053C^2 \tag{4}$$

Eq. (4) represents the combined usability deduced from the ANOVA model after analysis of data gathered from activities performed by participants.



Fig. 4. Normal probability.

Fig. 4 shows that all values of usability are normally plotted on or near the normal line. Its mean ANOVA model and RSM design are significant.

*E. Comparison between Theoretical and Actual Usability*

Fig. 5 represents the variations and comparison of perceptional usability, IT user usability and combined usability. According to IT users and experts, the usability is very close to the standard usability, combined usability is below the standard usability and users think that there is much need for improvements in the usability of MMAPPs.

*F. Analysis via Experimental Design II*

For this reason, 677 individuals concurred with the examination performed for twelve elements of mobile medical applications regarding point of interaction. The studies posed various inquiries about convenience assessment in settings with easiness, learnability, memorability, adequacy, productivity and satisfaction [36-37]. The study uses Google Forms for gathering their reactions because of the coronavirus pandemic,

likewise, lead Zoom meetings for direction about this review and top of the reactions. The study additionally guides through SMS, WhatsApp messages and calls to our IT specialists and investors about the survey on Google Forms [38]. Fig. 5 shows that improvement is required in all the attributes of usability as every group or run has usability value in the range of 50 to 64 for combined users, less than 60% for illiterate users and an average of nearly 75% for IT people. As per the above outcomes and conversation, there is a lot of progress expected to foster completeness and efficiency and reduce the complexity of MMAPPs. From the above conversation obviously, there is a hole among users, user perception level and application developers in setting with convenience which ought to be taken out by understanding the necessities and prerequisites of the users [39-40]. There is additionally an idea during advancement that the study might present the mode idea as designer presented in another application programming like Master Mode for doctors, User Mode for illiterate users and Well-disposed Mode for IT users.



Fig. 5. Comparison between IT users' perceptional and combined usability.

*G. Analytical Suggestions*

From the above conversation and information assembled, the study ought to likewise give all elements of medical in medical mobile applications for the fulfillment and viability of MMAPPs and conduct pieces of training for MMAPP's users. As indicated by the specialist's assessment and information accumulated from users, there are three segments, where upgradation is required which are complexity, efficiency and completeness.

*H. Limitations*

- The study took samples from one country, it may be extended worldwide.

- This study uses nine attributes of usability, which may increase to get more precision level.

## VI. CONCLUSION

Healthcare is evolving as the industry undergoes significant change. It is simple and advantageous for patients to adopt a healthy lifestyle by using mobile medical apps. The study worked on perceptional usability with the newly introduced UCD model with different attributes and checked their effect on usability through RSM designed model. Then analyzed and validated the model by ANOVA. On the other hand, traditional usability attributes are checked in a new form through the second RSM design model for IT users and combined user usability. It also checked the effects of attributes on perceptional and IT users and combined user usability as shown in Fig. 2 and Fig. 3. In the last compared these usability results of different groups of users according to the given design model. The deduced results of this review show that it is vital to think about patient fulfillment and confidence in MMAPPs for the future improvement of versatile medical application interfaces. Less than 50% of smartphone users utilize mobile medical applications (other than experts) to perform medical-related tasks for maintaining a healthy lifestyle. For this purpose, improve the awareness, make it much more interesting and enhance satisfaction, completeness, efficiency and easiness level of MMAPPs for users. The relationship between the average usage of medical applications and user health is statistically significant. A large number of participants agreed that medical applications can help to improve their health and as well as a healthy environment. It is necessary to limit the gap between patients/users and specialists for the improvement of MMAPPs. The last one is the expansion of all medical elements to further develop user fulfillment and user accommodation. There is no need to make many secure and complex MMAPPs like banking apps and security apps etc. At the end of the analysis, each design model produced a mathematical equation to evaluate its usability.

## REFERENCES

[1] M. Tanveer, H. Kaur, G. Thomas, H. Mahmood, M. Paruthi et al., "Mobile phone buying decisions among young adults: An empirical study of influencing factors," Sustainability, vol. 13, no. 19, pp. 19, 2021.

[2] J. Iqbal, N. Qureshi, M. A. Ashraf, S. F. Rasool and M. Z. Asghar, "The effect of emotional intelligence and academic social networking sites on academic performance during the COVID-19 pandemic," Psychol. res. behav. manag., vol. 14, no. 19, pp. 905–920, 2021.

[3] K. Hamid, M. W. Iqbal, H. A. B. Muhammad, Z. Fuzail, Z. T. Ghafoor et al., "Usability evaluation of mobile banking applications in digital business as emerging economy," International Journal of Computer Science and Network Security, vol. 22, no. 1, pp. 250–260, 2022.

[4] J. Businge, M. Openja, D. Kavaler, E. Bainomugisha, F. Khomh et al., "Studying android app popularity by cross-linking GitHub and google play store," in 2019 IEEE 26th international conference on software analysis, evolution and reengineering (SANER), Hangzou, HZ, China, pp. 287–297, 2019.

[5] A. Hussain, H. I. Abubakar and N. B. Hashim, "Evaluating mobile banking application: usability dimensions and measurements," in proceedings of the 6th international conference on information technology and multimedia, Putrajaya, PJ, Malaysia, pp. 136–140, 2014.

[6] F. Zahra, A. Hussain and H. Mohd, "Usability evaluation of mobile applications; where do we stand?," AIP conf. proc., vol. 1891, no. 1, pp. 020056, 2017.

[7] P. Jesilow, H. N. Pontell, and G. Geis, *Prescription for Profit: How Doctors Defraud Medicaid*. University of California Press, 2023.

[8] K. Hamid, M. waseem Iqbal, H. Muhammad, Z. Fuzail, and † Z., "ANOVA based usability evaluation of kid's mobile apps empowered

[9] J. Park and M. Zahabi, "A novel approach for usability evaluation of mobile applications," Proc. hum. factors ergon. soc. annu. meet., vol. 65, no. 1, pp. 437–441, 2021.

[10] J. Businge, M. Openja, D. Kavaler, E. Bainomugisha, F. Khomh et al., "Studying android app popularity by cross-linking GitHub and google play store," in 2019 IEEE 26th international conference on software analysis, evolution and reengineering (SANER), Hangzou, HZ, China, pp. 287–297, 2019.

[11] H. Cho, P. Y. Yen, D. Dowding, J. A. Merrill and R. Schnall, "A multi-level usability evaluation of mobile health applications: A case study," J. biomed. inform., vol. 86, no. 1, pp. 79–89, 2018.

[12] N. L. Hashim and A. J. Isse, "Usability evaluation metrics of tourism mobile applications," J. softw. eng. appl., vol. 12, no. 7, pp. 7, 2019.

[13] Sunardi, G. F. P. Desak, and Gintoro, "List of most usability evaluation in mobile application: a systematic literature review," in 2020 International Conference on Information Management and Technology (ICIMTech), Bandung, BD, Indonesia, pp. 283–287, 2020.

[14] H. I. Abubakar, N. L. Hashim and A. Hussain, "Usability evaluation model for mobile banking applications interface: model evaluation process using experts' panel," Journal of telecommunication, electronics and computer engineering, vol. 8, no. 10, pp. 53–57, 2016.

[15] K. Hamid, M. W. Iqbal, Z. Nazir, H. A. B. Muhammad, and Z. Fuzail, "Usability empowered by user's adaptive features in smart phones: the RSM approach,," Journal of Tianjin University, vol. 55, no. 7, pp. 285–304, 2022.

[16] U. E. M. Shah and T. K. Chiew, "A systematic literature review of the design approach and usability evaluation of the pain management mobile applications," Symmetry, vol. 11, no. 3, pp. 3–15, 2019.

[17] M. Z. Ansaar, J. Hussain, J. Bang, S. Lee, K. Y. Shin, and K. Young Woo, "The mHealth Applications Usability Evaluation Review," in 2020 International Conference on Information Networking (ICOIN), Jan. 2020, pp. 70–73. doi: 10.1109/ICOIN48656.2020.9016509.

[18] H. Cho, P.-Y. Yen, D. Dowding, J. A. Merrill, and R. Schnall, "A multi-level usability evaluation of mobile health applications: A case study," Journal of Biomedical Informatics, vol. 86, pp. 79–89, Oct. 2018, doi: 10.1016/j.jbi.2018.08.012.

[19] P. M. A. B. Estrela, R. D. O. Albuquerque, D. M. Amaral, W. F. Giozza and R. T. D. S. Júnior, "A framework for continuous authentication based on touch dynamics biometrics for mobile banking applications," Sensors, vol. 21, no. 12, pp. 12–24, 2021.

[20] J. Park and M. Zahabi, "A novel approach for usability evaluation of mobile applications," Proc. hum. factors ergon. soc. annu. meet., vol. 65, no. 1, pp. 437–441, 2021.

[21] J. Orlovska, C. Wickman and R. Söderberg, "Big data analysis as a new approach for usability attributes evaluation of user interfaces: an automotive industry context," in DS 92: Proceedings of the DESIGN 2018 15th International Design Conference, Edinburg, EB, Scotland, pp. 1651–1662, 2018.

[22] N. M. Zaharakis, M. J. Mason and C. Berkel, "Responsiveness to mHealth intervention for cannabis use in young adults predicts improved outcomes," Prev. sci., vol. 23, no. 4, pp. 630–635, 2022.

[23] A. Riegler and C. Holzmann, "Measuring visual user interface complexity of mobile applications with metrics," Interact. comput., vol. 30, no. 3, pp. 207–223, 2018.

[24] K. Hamid, H. Muhammad, M. waseem Iqbal, A. Nazir, shazab, and H. Moneeza, "ML-Based Meta Model Evaluation of Mobile Apps Empowered Usability of Disables," *Tianjin Daxue Xuebao Ziran Kexue Yu Gongcheng Jishu BanJournal Tianjin Univ. Sci. Technol.*, vol. 56, pp. 50–68, Jan. 2023.

[25] M. W. Iqbal, N. A. Ch, S. K. Shahzad, M. R. Naqvi, B. A. Khan et al., "User context ontology for adaptive mobile-phone interfaces," IEEE access, vol. 9, no. 1, pp. 96751–96762, 2021.

[26] K. Hamid, H. Muhammad, M. waseem Iqbal, S. Bukhari, A. Nazir, and S. Bhatti, "ML-Based Usability Evaluation of Educational Mobile Apps for Grown-Ups and Adults," *Jilin Daxue Xuebao GongxuebanJournal*

learning process," Qingdao Daxue Xuebao(Gongcheng Jishuban)/Journal of Qingdao University (Engineering and Technology Edition), vol. 41, no. 6, pp. 142–169, 2022.

*Jilin Univ. Eng. Technol. Ed.*, vol. 41, pp. 352–370, Dec. 2022, doi: 10.17605/OSF.IO/YJ2E5.

[27] K. Hamid, M. waseem Iqbal, Z. Nazir, H. Muhammad, and Z. Fuzail, "Usability Empowered by User's Adaptive Features in Smart Phones: The RSM Approach," *Tianjin Daxue Xuebao Ziran Kexue Yu Gongcheng Jishu BanJournal Tianjin Univ. Sci. Technol.*, vol. 55, pp. 285–304, Jul. 2022, doi: 10.17605/OSF.IO/6RUZ5.

[28] M. Aqeel *et al.*, "Response Surface Methodology-Based Usability Evaluation of Apps for Visually Impaired Persons," vol. 42, pp. 532–545, Mar. 2023, doi: 10.17605/OSF.IO/7G29Z.

[29] N. C. Rust and V. Mehrpour, "Understanding image memorability," Trends in cognitive sciences, vol. 24, no. 7, pp. 557–568, 2022.

[30] M. Iqbal, N. Ahmad and S. K. Shahzad, "Usability evaluation of adaptive features in smartphones," procedia computer science, vol. 112, no. 1, pp. 2185–2194, 2017.

[31] K. A. Sespiani and N. F. Ernungtyas, "Connecting elderly and digital devices: a literature review of user interface studies for indonesian elders," J. soc. media, vol. 6, no. 1, pp. 139-156, 2022.

[32] H. I. Abubakar, N. L. Hashim and A. Hussain, "Usability evaluation model for mobile banking applications interface: model evaluation process using experts' panel," Journal of telecommunication, electronics and computer engineering, vol. 8, no. 10, pp. 53–57, 2016.

[33] M. Iqbal, N. Ahmad and S. K. Shahzad, "Usability evaluation of adaptive features in smartphones," Procedia comput. sci., vol. 112, no. 1, pp. 2185–2194, 2017.

[34] A. Hodrien, and T. P. Fernando, "A review of post-study and post-task subjective questionnaires to guide assessment of system usability," Journal of usability studies, vol. 16, no. 1, pp. 203–232, 2021.

[35] E. Boeren and T. Í. Berrozpe, "Unpacking PIAAC's cognitive skills measurements through engagement with bloom's taxonomy," Studies in educational evaluation, vol. 73, no. 1, pp. 101–151, 2022.

[36] F. K. Mazumder, "Usability guidelines for usable user interface," International Journal of Research Engineering and Technology, vol. 03, no. 09, pp. 79–82, 2014.

[37] V. J. Aski, , V. S. Dhaka, , S. Kumar, , S. Verma and D. B. Rawat, "Advances on networked ehealth information access and sharing: status, challenges and prospects," Computer networks, vol. 204, no. 1, pp. 108687, 2022..

[38] A. W. Siyal, D. Donghong, W. A. Umrani, S. Siyal and S. Bhand, "Predicting mobile banking acceptance and loyalty in chinese bank customers," SAGE open, vol. 9, no. 2, pp. 2158244019844084, 2019.

[39] L. Tao and M. Zhang, "Understanding an online classroom system: design and implementation based on a model blending pedagogy and HCI," IEEE transactions hum.-mach. syst., vol. 43, no. 5, pp. 465–478, 2013.

[40] J. Tang, "Discussion on health service system of mobile medical institutions based on internet of things and cloud computing." Journal of Healthcare Engineering, vol. 2022, no. 5235349, pp. 12–25, 2022.

# Development of a Framework for Predicting Students' Academic Performance in STEM Education using Machine Learning Methods

Rustam Abdrakhmanov[1], Ainur Zhaxanova[2], Malika Karatayeva[3],
Gulzhan Zholaushievna Niyazova[4], Kamalbek Berkimbayev[5], Assyl Tuimebayev[6]

International University of Tourism and Hospitality, Turkistan, Kazakhstan[1]
South Kazakhstan University named after M.Auezov, Shymkent, Kazakhstan[2]
M. Auezov South Kazakhstan University, Shymkent, Kazakhstan[3]
Khoja Akhmet Yassawi International Kazakh-Turkish University, Turkistan, Kazakhstan[4, 5]
Boston University, Boston, USA[6]

*Abstract*—In the continuously evolving educational landscape, the prediction of students' academic performance in STEM (Science, Technology, Engineering, Mathematics) disciplines stands as a paramount component for educational stakeholders aiming at enhancing learning methodologies and outcomes. This research paper delves into a sophisticated analysis, employing Machine Learning (ML) algorithms to predict students' achievements, focusing explicitly on the multifaceted realm of STEM education. By harnessing a robust dataset drawn from diverse educational backgrounds, incorporating myriad factors such as historical academic data, socioeconomic demographics, and individual learning interactions, the study innovates by transcending traditional prediction parameters. The research meticulously evaluates several machine learning models, juxtaposing their efficacies through rigorous methodologies, including Random Forest, Support Vector Machines, and Neural Networks, subsequently advocating for an ensemble approach to bolster prediction accuracy. Critical insights reveal that customized learning pathways, preemptive identification of at-risk candidates, and the nuanced understanding of contributing influencers are significantly enhanced through the ML framework, offering a transformative lens for academic strategies. Furthermore, the paper confronts the ethical quandaries and challenges of data privacy emerging in the wake of advanced analytics in education, proposing a holistic guideline for stakeholders. This exploration not only underscores the potential of machine learning in revolutionizing predictive strategies in STEM education but also advocates for continuous model optimization, embracing a symbiotic integration between pedagogical methodologies and technological advancements, thereby redefining the trajectories of educational paradigms.

*Keywords—Load balancing; machine learning; server; classification; software*

## I. INTRODUCTION

In the current epoch of technological ubiquity, the domain of education, particularly Science, Technology, Engineering, and Mathematics (STEM) education, has encountered transformative shifts. The imperative to mold proficient future professionals capable of navigating complex technological terrains and scientific quandaries has never been more pressing [1]. Yet, the chasm between educational methodologies and individual student performance continues to challenge educators and policy-makers alike, necessitating innovative approaches to bridge this gap. Central to this innovation is the utilization of machine learning (ML) [2] in comprehending and predicting student performance, a research niche that has burgeoned in significance due to its profound implications on educational strategies [3].

Historically, educational outcomes were often predicated on conventional metrics—standardized testing, classroom participation, and rudimentary performance tracking methodologies [3]. These linear models, although somewhat informative, hardly capture the labyrinth of individual student experiences, inherent talents, cognitive styles, and external factors impacting academic performance. The intricacies of learning are often obfuscated by the homogeneity of traditional assessment tools, which are ill-equipped to forecast academic outcomes with substantial reliability [4]. The need for personalized education, a clarion call in contemporary pedagogical circles, further exacerbates this issue, as traditional educational models are systemically inept at accommodating the heterogeneity of student populations [5].

Emerging from this backdrop is the promise of machine learning, a subset of artificial intelligence (AI) characterized by its capacity for pattern recognition, learning from data, and making predictions [6]. When applied within the educational sphere, ML bears the revolutionary potential to distill vast, nebulous datasets into actionable insights regarding student performance. This process is not without its complexities, as it necessitates a delicate alchemy of algorithmic selection, hyperparameter tuning, and feature engineering, demanding rigorous scrutiny to ensure both ethical and practical efficacy [7].

Literature in the realm of machine learning applications in education is replete with instances of predictive analytics being employed for student data. Studies range from early identification of students at risk of academic failure to nuanced understandings of how socio-economic factors correlate with educational outcomes [8]. Specifically, within STEM disciplines, where abstract concepts and cumulative learning

are pivotal, the predictive power of ML can aid in identifying learning hurdles and pedagogical inefficiencies [9].

However, despite its burgeoning presence, the integration of machine learning into educational predictive models is fraught with challenges. The ethical implications of data privacy, security, and the potential for bias in algorithmic determinations present formidable hurdles [10]. Each of these aspects requires careful consideration to maintain the integrity of educational institutions while harnessing the capabilities of advanced technology. Moreover, there is the omnipresent challenge of interpretability, as the decision-making processes of complex models often constitute a "black box," making it difficult for educators and stakeholders to trust and ethically utilize the predictions [11].

This research, therefore, is anchored in the critical evaluation of various machine learning models in predicting student performance in STEM subjects. The choice of models, including but not limited to, Random Forest, Support Vector Machines, and Neural Networks, represents the spectrum of algorithms from simple interpretable models to complex, high-dimensional ones, each with unique strengths and predictive accuracies [12]. Furthermore, the study leverages an ensemble learning approach, conjectured to enhance the robustness and reliability of predictions through the aggregation of multiple models [13].

The nuance of this research resides in its holistic approach, not just considering academic datasets but also integrating comprehensive student data. This encompasses demographic information, previous academic achievements, engagement levels, and even socio-economic indicators, acknowledging the multifactorial nature of educational success [14]. By doing so, the research transcends myopic academic predictions, offering instead a panoramic view of student performance influencers. This approach is pivotal, recognizing that contemporary students navigate a milieu replete with both academic and non-academic challenges, ranging from mental health pressures to the digital distractions endemic in modern society [15].

In synthesizing these elements, this study contributes to the academic dialogue in several ways. Firstly, it provides an empirical evaluation of machine learning models in the context of education, a field where such advanced technology applications remain under-explored. Secondly, it addresses the ethical and practical challenges intrinsic to the domain, offering pathways for stakeholders to leverage insights responsibly. Finally, by focusing on STEM education—a critical driver of future innovation and economic growth—the research underscores the need for educational systems to evolve in tandem with broader societal advancements, ensuring that student success is not left to antiquation in this brave new world [16].

In essence, this paper seeks to navigate the confluence of technology and education, providing insights that could potentially reshape the predictive paradigms in the educational sector, particularly within STEM disciplines. Through rigorous analysis, ethical considerations, and practical applications, the study stands as a beacon, guiding the way towards a more informed, equitable, and effective educational landscape.

## II. RELATED WORKS

The integration of machine learning (ML) strategies in educational settings, particularly in STEM (Science, Technology, Engineering, Mathematics) education, has garnered considerable attention in academic circles, underpinning the necessity to comprehend its precedents within scholarly literature. The endeavor to harness ML's predictive power to forecast student academic performance intersects various research domains, necessitating an interdisciplinary approach to fully understand its scope, potential, and limitations [17].

### A. Early Interventions and Predictive Analytics

The genesis of employing predictive analytics in education was primarily focused on identifying at-risk students to facilitate early interventions. Studies by [18] and [19] showcased the utility of traditional statistical methods to forecast academic struggles, primarily utilizing historical and continuous assessment data. However, these models often suffered from oversimplification, failing to capture the multifaceted nature of academic performance. The advent of machine learning offered a nuanced approach, enabling the consideration of a broader array of variables and revealing hidden patterns within complex datasets [20]. For instance, [21] exploited decision trees to flag students needing additional support, demonstrating superior accuracy over classical methods.

### B. Machine Learning Models in Academic Settings

Diverse ML models have been tested within educational contexts, each offering distinct advantages. Research by [22] emphasized the effectiveness of Support Vector Machines (SVM) in handling high-dimensional spaces, particularly useful in analyzing voluminous student data. In contrast, studies like [23] touted Random Forest's ability to provide insight into feature importance, thereby understanding influential factors affecting student performance. Neural Networks, known for their prowess in capturing non-linear relationships, were explored by [24], underscoring their sensitivity to nuanced interactions amongst data, albeit at the expense of interpretability.

### C. Ensemble Methods and Hybrid Models

Given the heterogeneous nature of educational data, recent studies have advocated for ensemble and hybrid models. [25] demonstrated that combining predictions from multiple algorithms enhanced overall accuracy, compensating for individual model weaknesses. Similarly, [26] successfully implemented a hybrid model incorporating both neural networks and decision trees, exploiting the strengths of both non-linearity and interpretability. These methodologies signify a shift toward more robust predictive systems, though they necessitate careful construction and validation.

### D. Behavioral and Engagement Metrics

Beyond academic results, behavioral and engagement metrics have surfaced as vital predictors. [27] integrated data from online learning platforms, highlighting that digital engagement levels were indicative of academic outcomes. This was corroborated by [28], who found that behavioral patterns, such as time spent on tasks and participation in virtual forums,

were salient in understanding academic performance. These insights underscore the importance of a holistic data approach, merging academic, behavioral, and engagement metrics.

### E. Socio-economic and External Factors

Acknowledging the impact of external factors, several researchers have broadened the data scope to include socio-economic factors. The research in [29] affirmed that socio-economic status significantly correlated with academic achievement, while [30] extended this by showing that even when controlling for this, other factors, including parental involvement and peer influence, played non-trivial roles. The study in [31] further incorporated these into a comprehensive ML model, illustrating the enhanced predictive power when acknowledging the multifactorial nature of education.

### F. Ethical Considerations and Bias Mitigation

The ethical dimensions of ML in education, especially concerning data privacy and algorithmic bias, have provoked intense scholarly discourse. The study in [32] critically analyzed ethical quandaries, stressing the need for transparency, consent, and privacy safeguards. The research in [33] explored the prevalence of biases, showing that unexamined, algorithms might perpetuate existing inequalities, necessitating rigorous bias mitigation protocols. The responsibility of ethical algorithm deployment is echoed throughout literature, demanding a balance between technological advancement and moral obligations [34].

### G. Interpretability and Decision Transparency

The "black box" nature of certain ML models presents substantial challenges in educational settings, where stakeholders require transparency to trust and act upon predictions. Research by [35] proposed methodologies for enhancing the interpretability of complex models, while [36] discussed the trade-offs between accuracy and interpretability, suggesting that simpler models might sometimes serve educational needs better due to their transparency.

### H. Customized Learning Pathways

Tailoring education to individual needs is another frontier. [37] demonstrated that ML could help create customized learning pathways, thereby improving engagement and comprehension. This personalization aspect, especially in STEM subjects that often suffer from high drop-out rates, can potentially revolutionize educational methodologies [38].

### I. Challenges and Future Directions

Despite its promise, the integration of ML in education isn't without challenges. The study in [39] outlined issues ranging from data quality, privacy concerns, and the need for interdisciplinary collaboration between educators and data scientists. The literature strongly advocates for continued research, particularly iterative model refinement and the exploration of innovative data sources to enrich predictive capabilities.

In summary, the existing body of work confirms the transformative potential of machine learning in predicting academic performance in STEM education. It reflects a trajectory from simplistic predictive models towards more sophisticated, comprehensive, and ethically considerate ML

applications. This literature tapestry provides a foundation upon which the current research is built, aiming to contribute novel insights by harnessing the potency of ML to navigate the complex, dynamic landscape of educational predictors and outcomes.

## III. MATERIALS AND METHODS

The methodological framework guiding this research is visually represented in Fig. 1, elucidating a comprehensive five-stage process integral to the operationalization of this study. Initially, the process commences with the meticulous aggregation of relevant datasets, sourced extensively from institutional databases, ensuring a rich compendium of variables reflective of the educational milieu. Subsequently, the study introduces a sophisticated application of natural language processing (NLP) techniques, aimed at dissecting and quantifying classroom dialogic interactions, a step that underscores the significance of linguistic dynamics in educational settings.



Fig. 1. Data collection and preparation.

Progressing beyond raw data compilation, the research methodologically embraces rigorous statistical methodologies to scrutinize dialogue-based indicators, thereby quantifying abstract elements of classroom discourse. This transformative approach facilitates a nuanced understanding of pedagogical dynamics, often overlooked in traditional analysis paradigms. In the ensuing phase, the study leverages state-of-the-art deep learning algorithms, architecting a predictive model for academic performance that is both robust and sensitive to the multifarious factors influencing educational outcomes.

The final step epitomizes the study's commitment to transparency and applicability through the adoption of an interpretable artificial intelligence (AI) model. This phase is dedicated to the explication of critical predictors within the established predictive model, a crucial aspect that not only enhances the trustworthiness of AI interventions but also empowers stakeholders with actionable insights. The ensuing sections are committed to an in-depth exposition of each pivotal stage, shedding light on the intricate methodologies that constitute the backbone of this research endeavor, thereby

reinforcing its academic rigor and practical relevance in the educational echelon.

### A. Dataset

Data for this study were meticulously sourced from virtual classrooms within a prominent online educational framework in Kazakhstan. These live classrooms, distinctive in their interactive nature, allow students to exchange messages visible to their peers, thereby fostering a dynamic communicative environment through an integrated chat room feature. The research harnessed transcripts of these real-time educational dialogues across various subjects and academic levels, specifically focusing on the platform's courses for grades K-6.

Reflecting on the outcomes of the 2020 spring semester, the platform recorded an enrollment of approximately 30,581 students within the K-6 category. Of these, around 8,158 students were engaged in 2,545 courses unrelated to STEM, while a notably larger cohort of 22,423 students immersed themselves in 3,797 STEM-oriented courses. This academic engagement led to the creation of approximately 654,954 and 1,690,549 interactive textual exchanges for non-STEM and STEM courses, respectively. Table I provides a detailed breakdown of the dialogue texts within the live classroom chat environments, with 'M' denoting the mean value.

In assessing academic performance, this study adopted a comparative analysis of students' rankings, focusing on discrepancies between their initial (pretest) and final (posttest) standings. The upper echelon of academic achievers, represented by the top 20%, was classified as the high-performance group, whereas the lowermost 20% of the spectrum was identified as the low-performance group. For subsequent analytical procedures, the study incorporated data from approximately 4,776 students—around 2,459 from the low-performance segment and 2,317 from the high-performance tier—in non-STEM courses. Simultaneously, the STEM counterparts comprised a more substantial ensemble of

roughly 13,659 students, segmented into 7,711 underachievers and 5,948 high achievers. This strategic dichotomy in performance assessment is instrumental in facilitating a nuanced understanding of educational dynamics within the virtual classroom scenario.

### B. Applying Machine Learning

To facilitate the automated identification of emotional articulations and the categorization of interactive modalities within classroom dialogues, we have dedicated efforts toward the development and training of two distinct models of text classification. This intricate process, essential for comprehending the underpinnings of communicative exchanges in educational settings, is graphically synthesized in the flow diagram presented as Fig. 2. This visual representation underscores the systematic approach adopted for this phase of the research, highlighting the advanced computational techniques employed to analyze textual data within the pedagogical discourse.

In this study, we established a nuanced criterion for the categorization of emotional tenor and interaction modalities, substantiated with specific textual instances derived directly from classroom dialogues. Emotional expressions within the communicative exchanges are bifurcated into two primary affective states: positive and negative emotions. Dialogues permeated with sentiments of joy, elation, or exhilaration is classified under the umbrella of positive emotional discourse. Conversely, student interactions manifesting tones of melancholy, disinterest, or irritation are categorized as expressions of negative emotion. This dichotomous approach to emotional categorization provides a streamlined yet profound understanding of the affective landscape of classroom interactions. Representative samples of these categorized emotional states, extracted verbatim from the dialogic exchanges, are systematically presented in Table II, offering tangible insights into the emotional substrates that underpin student communication within academic settings.

TABLE I.  DIALOG TESTS IN CLASSROOM

| Course | Subject | Grade | Classes | M(SD) | | Number of interactive course |
|---|---|---|---|---|---|---|
| | | | | Number of students | Number of interactive texts | |
| Non-STEM subject | English | 1 | 598 | 1.97 | 98 | 93 467 |
| | | 2 | 867 | 5.39 | 61 | 281 679 |
| | | 3 | 2387 | 2.8 | 79 | 564 782 |
| | | 4 | 1647 | 3.7 | 90 | 483164 |
| | | 5 | 1983 | 3.5 | 86 | 582 264 |
| | | 6 | 1976 | 3.6 | 78 | 564 778 |
| STEM subject | Math | 1 | 2729 | 7.9 | 6159 | 1 067 899 |
| | | 2 | 4318 | 6.8 | 80 | 1 647 896 |
| | | 3 | 3016 | 5.6 | 99 | 1 302 445 |
| | | 4 | 1725 | 3.8 | 111 | 631 448 |
| | | 5 | 1973 | 5.8 | 123 | 1 305 866 |
| | | 6 | 1609 | 5.5 | 137 | 1 018 886 |

Fig. 2. Dialog classification schema.

TABLE II. EXAMPLES OF STUDENT EXPRESSIONS

| Dimension | First-level | Second level | Sample |
|---|---|---|---|
| Expression | Positive | - | I like it very much |
| | Negative | - | It was boring to me; It is not interesting |
| Interactive | Cognition | Asking questions and answering | What is the meaning of … |

### C. Academic Performance Prediction

In the context of this research, features integral to the construction of predictive models are bifurcated into two distinct categories. The initial category encompasses the 'pretest rank,' signifying students' foundational academic competencies prior to their engagement in specific classroom sessions. The subsequent category is more intricate, involving features meticulously extrapolated from the text of classroom dialogues. Collectively, these categories amalgamate into a robust set of 48 distinctive features, instrumental in the subsequent phases of predictive model construction.

The study employs a comparative analysis approach, rigorously evaluating three sophisticated classification algorithms pivotal for predicting academic performance. These encompass a Convolutional Neural Network (CNN)

methodology as illustrated in Fig. 3. The model is subjected to an in-depth assessment based on three critical evaluation metrics: recall, precision, and accuracy. This triadic evaluative framework provides a holistic view of each algorithm's performance, thereby informing the selection of the most efficacious predictive model.

Upon the empirical evaluation of these algorithms, the study advances to synthesize an interpretable model, enhancing the applicability and user comprehension of the results. Notably, the implementation phase of this research utilizes the TensorFlow deep learning framework, a decision substantiated by its proven efficacy and robustness in handling complex predictive tasks. This strategic methodological orchestration not only underscores the rigor of the study but also enhances the reliability and validity of the predictive outcomes within the academic performance landscape.



Fig. 3. Deep learning for predicting academic performance of students.

## IV. EXPERIMENTAL RESULTS

### A. Evaluation of Emotional Expression of Students

This study engages with the SHAP (SHapley Additive exPlanations) methodology, an advanced technique within the realm of interpretable artificial intelligence, to critically analyze the contributory features inherent in the academic performance prediction model. The consequential insights derived from this rigorous analysis are graphically elucidated in Fig. 3. Concurrently, an intriguing observation emerges from the C1 cohort, exhibiting a marginal enhancement in predictive accuracy relative to the established baseline, which is preliminarily set at 50%. This nuanced increment, albeit minimal, signals a critical inference: the interactive dynamics encapsulated within the online classroom environment exert a relatively insubstantial influence on the academic trajectories associated with non-STEM coursework. This revelation underscores the necessity for a differential pedagogical approach, potentially customized to the distinct educational exigencies of STEM and non-STEM curricula.

Intricately woven into this analysis are six pivotal variables, each derived from a comprehensive aggregation of the absolute values of corresponding interactive or emotional metrics within a specific interactive phase. For instance, the variable 'summary_interaction' is computed by summing the absolute SHAP values of various interactive categories during the summary stage, represented formulaically as: summary_interaction = $|ics| + |ims| + |ios| + |ccs| + |cms| + |cos|$. Analogously, 'summary_emotion' encapsulates the emotional undertones of the summary phase, calculated as: summary_emotion = $|ips| + |cps| + |ins| + |cns|$.

These computations underscore the nuanced complexity and the multifaceted nature of interactive and emotional dynamics within the learning environment. By leveraging SHAP values, the study provides an in-depth, interpretable analysis, highlighting the often-overlooked subtleties that significantly influence students' academic trajectories in non-STEM disciplines. This meticulous approach not only enhances the comprehensibility of predictive analytics but also informs educational strategies by pinpointing specific areas of student-teacher interaction that require pedagogical attention.

### B. CNN in Academic Performance Prediction

In the methodological framework of our research, we meticulously partitioned the dataset, allocating 70% to training purposes, while the residual 30% was designated as the test set, concurrently serving as the validation data. It is pertinent to clarify that within the context of this study, the terminologies 'test loss' and 'test accuracy' are utilized interchangeably with 'validation loss' and 'validation accuracy,' respectively. Our evaluative metric of paramount importance was accuracy, a choice that steered the subsequent analytical processes, including the imperative exercise of parameter optimization, to elicit the most robust and reliable outcomes.

One critical parameter that demanded our focused attention was the learning rate, recognized for its decisive role in the training phase of machine learning models. Typically, a diminutive positive number confined within the spectrum of 0 to 1, the learning rate orchestrates the velocity at which a

model acclimatizes to a given problem. Its optimal calibration is crucial; an excessively accelerated learning rate might precipitate a premature convergence, culminating in a suboptimal solution, whereas a rate set too sluggishly risks miring the process in stagnation.

Given these potential quandaries, our study was committed to identifying the most propitious learning rate. We embarked on empirical trials employing a gamut of learning rates, meticulously observing their impacts on model performance. Further augmenting our analysis, we crafted visual representations of the correlation between diverse learning rates and their corresponding training and test accuracies. These illustrative delineations, accessible in Fig. 4, not only enhance comprehensibility but also provide empirical substantiation for the optimal learning rate conducive to our model's most effective learning trajectory.

An insightful observation emerges from the graphical representation delineated in Fig. 4, wherein the learning curve exhibits a pronounced stagnation at elevated learning rates, specifically at 1 and 0.1. This phenomenon indicates an inhibited learning process, attributable to the model's inability to assimilate the training data effectively at such escalated rates. Conversely, the curves corresponding to reduced learning rates reveal a propensity for oscillation, a manifestation of inconsistent learning. Within the context of our experimental framework, empirical evidence converges on the learning rate of 0.001 as the optimal parameter, a conclusion corroborated by the enhanced performance metrics of the model discernible in Fig. 4.



Fig. 4. Training and test accuracy of the applied model in 100 learning epochs.

An insightful observation emerges from the graphical Furthermore, our research extended into an exploratory analysis involving varying numbers of epochs, intended to ascertain their influence on the model's accuracy. This phase entailed methodical observations of the trajectories of accuracy versus epochs, along with loss versus epochs, intensifying our understanding of the model's behavior over iterative learning cycles. The experimentation commenced with an epoch baseline of 20, progressively amplifying to discern the consequential impacts on model performance.

The empirical outcomes derived from this procedural iteration, particularly in terms of test loss and test accuracy, are of significant interest. Fig. 5 in the study play a crucial role in this context. These figures provide a visual representation of the behavior of accuracy and loss across various epoch parameters, respectively. By examining these figures, one can gain a deeper understanding of how the model's accuracy and loss metrics evolve throughout the training process.



Fig. 5. Training and test accuracy of the applied model in 30 learning epochs.

Together, these figures provide a comprehensive view of the model's learning consistency and its predictive performance. By analyzing these graphs, researchers can determine the most effective epoch configuration, balancing the need for sufficient training to achieve high accuracy without overfitting. This balance is critical for ensuring that the model remains generalizable and performs well on new, unseen data.

## V. DISCUSSION

This research embarked on an intricate journey to unravel the layers of complexities in predicting students' academic performance within the STEM education landscape, utilizing the prowess of machine learning algorithms. The findings illuminate several critical facets of educational psychology, pedagogical strategies, and the subtle nuances of student interactions and engagement, particularly in an online learning environment.

One of the cardinal revelations of this study was the pivotal role of interactive patterns and emotional expressions in shaping students' academic outcomes. Previous research confined itself to traditional performance indicators, often overlooking the rich tapestry of student interactions. Our study bridged this gap, echoing the findings of [40], which underscored the significance of emotional and psychological factors in academic performance. However, unlike [41] that generalized the impact of interactive patterns, our research unveiled a stage-wise influence, emphasizing that the timing of interactions is just as crucial as their nature.

Furthermore, the disparity in the influence of these interactive factors between STEM and non-STEM courses is particularly enlightening. Consistent with the observations of [42], our findings corroborate that STEM subjects, with their structured and logical framework, respond differently to emotional and interactive stimuli compared to non-STEM subjects. This nuanced understanding advocate for a more tailored approach in pedagogical strategies, as also suggested by [43-44], ensuring that educators can mold their teaching tactics according to the subject matter and the corresponding emotional and interactive dynamics.

The application of machine learning, especially deep learning algorithms, marked a paradigm shift in identifying and predicting successful learning patterns. While traditional statistical methods provided a surface-level understanding, the neural networks delved deeper into the data, much like the human brain, offering unprecedented insights into student performance predictors. This sophisticated approach, however, came with its own set of challenges, chiefly selecting the appropriate model and tuning the hyperparameters, as discussed in [45].

Our research determined the optimal learning rate, a finding that resonated with the work of [46], highlighting the delicate balance required in setting this parameter. Too high a rate, and the model overshoots the minimum point; too low, and the model succumbs to local minima or becomes computationally impractical. Similarly, the number of epochs represented a tug of war between underfitting and overfitting, a common conundrum in machine learning models as identified by [47]. Our study struck this balance adeptly, ensuring the model learned the underlying patterns without memorizing the data, a nuanced mastery over the art of 'learning to learn'.

Interestingly, the efficacy of the model was not universally uniform across different subject matters. While it showed remarkable precision in predicting STEM outcomes, its applicability in non-STEM subjects was limited, a phenomenon that could be attributed to the inherent subject differences. STEM subjects, often characterized by logical and structured learnings, lend themselves more readily to predictive analytics, unlike non-STEM subjects that are more abstract and open to interpretation, as noted by [48].

Another intriguing aspect was the identification of effective features in the predictive model using SHAP values, an area often shrouded in mystery in most machine learning applications due to their 'black box' nature. The interpretability introduced by SHAP values, as explored in [49], demystified the influential features, providing invaluable insights for educators. Knowing which factors are more indicative of a student's performance could revolutionize educational strategies, ensuring a more focused and student-centric approach.

However, despite these advancements, the limitation of data cannot be overlooked. While the study harnessed a wealth of data points, the quality of these data, especially concerning the emotional expressions, was heavily reliant on the accuracy of the text classification models. Future research could benefit from more sophisticated Natural Language Processing (NLP) tools, possibly incorporating contextual understanding to grasp the subtleties of human emotion and interaction, an enhancement suggested by [50].

Moreover, the scope of the dataset also posed a constraint. The research was circumscribed to a specific geographical region and educational level, limiting the generalizability of the findings. Subsequent studies could transcend these boundaries, encompassing a more diverse student population to authenticate the universality of the findings.

In conclusion, this research has paved a novel pathway in understanding and predicting students' academic performance, intertwining the realms of machine learning, educational strategies, and psychological underpinnings. The insights gleaned hold profound implications for educators, policy-makers, and curriculum designers, advocating for a more holistic, student-oriented approach in the educational odyssey. However, the journey does not end here. With the continual evolution of machine learning and the ever-changing educational landscape, future research beckons, promising even deeper insights and more personalized educational experiences.

## VI. CONCLUSION

The journey through this research, from conceptual frameworks to analytical discussions, reflects a profound exploration of integrating machine learning into Software-Defined Networking (SDN) to enhance load balancing. As we draw conclusions, it's imperative to encapsulate the essence of our findings and their implications for future scientific inquiry and practical application in the networking sphere.

This study marked a significant advancement by demonstrating that machine learning algorithms could revolutionize the way network resources are managed, optimizing the distribution of data loads across various pathways. By employing sophisticated algorithms, we unveiled the potential to predict network congestions, dynamically adjust to traffic changes, and improve overall efficiency and user experience. This paradigm shift from traditional methods accentuates a move towards more autonomous, self-sufficient systems capable of sophisticated decision-making processes, essential in the burgeoning era of digital transformation and the Internet of Things (IoT).

However, the research also highlighted critical challenges and limitations, from the complexities of algorithm training and data security concerns to the practical applicability of the proposed model outside simulated environments. These challenges are not terminuses but instead signposts indicating areas requiring further exploration, refinement, and innovation.

Looking forward, the implications of this research are both broad and profound. They suggest an imminent need for robust, real-world testing and the potential for interdisciplinary approaches that could further enrich these technological advances. The prospects of enhanced security measures, scalability considerations, and user-centric adaptations also present exciting, necessary trajectories.

In conclusion, this study does not represent an end but a beginning. It serves as a catalyst for continued exploration and dialogue in the realms of machine learning, networking, and beyond. The confluence of these fields holds significant promise for creating more resilient, efficient, and intelligent networks, poised to support the ever-evolving demands of future digital landscapes.

## REFERENCES

[1] Hamdan, M., Hassan, E., Abdelaziz, A., Elhigazi, A., Mohammed, B., Khan, S., ... & Marsono, M. N. (2021). A comprehensive survey of load balancing techniques in software-defined network. Journal of Network and Computer Applications, 174, 102856.

[2] Omarov, B., Omarov, B., Shekerbekova, S., Gusmanova, F., Oshanova, N., Sarbasova, A., ... & Sultan, D. (2019). Applying face recognition in video surveillance security systems. In Software Technology: Methods and Tools: 51st International Conference, TOOLS 2019, Innopolis, Russia, October 15–17, 2019, Proceedings 51 (pp. 271-280). Springer International Publishing.

[3] Nayak, R. P., Sethi, S., Bhoi, S. K., Sahoo, K. S., & Nayyar, A. (2023). Ml-mds: Machine learning based misbehavior detection system for cognitive software-defined multimedia vanets (csdmv) in smart cities. Multimedia Tools and Applications, 82(3), 3931-3951.

[4] Muhammad, T. (2022). A Comprehensive Study on Software-Defined Load Balancers: Architectural Flexibility & Application Service Delivery in On-Premises Ecosystems. International Journal of Computer Science and Technology, 6(1), 1-24.

[5] Rahman, A., Islam, J., Kundu, D., Karim, R., Rahman, Z., Band, S. S., ... & Kumar, N. (2023). Impacts of blockchain in software-defined Internet of Things ecosystem with Network Function Virtualization for smart applications: Present perspectives and future directions. International Journal of Communication Systems, e5429.

[6] Murugesan, G., Ahmed, T. I., Shabaz, M., Bhola, J., Omarov, B., Swaminathan, R., ... & Sumi, S. A. (2022). Assessment of mental workload by visual motor activity among control group and patient suffering from depressive disorder. Computational Intelligence and Neuroscience, 2022.

[7] Jurado-Lasso, F. F., Marchegiani, L., Jurado, J. F., Abu-Mahfouz, A. M., & Fafoutis, X. (2022). A survey on machine learning software-defined wireless sensor networks (ml-SDWSNS): Current status and major challenges. IEEE Access, 10, 23560-23592.

[8] Wu, J., Dong, M., Ota, K., Li, J., & Yang, W. (2020). Application-aware consensus management for software-defined intelligent blockchain in IoT. IEEE Network, 34(1), 69-75.

[9] Yazdinejad, A., Parizi, R. M., Dehghantanha, A., & Choo, K. K. R. (2020). P4-to-blockchain: A secure blockchain-enabled packet parser for software defined networking. Computers & Security, 88, 101629.

[10] Karakus, M., Guler, E., & Uludag, S. (2021). Qoschain: Provisioning inter-as qos in software-defined networks with blockchain. IEEE Transactions on Network and Service Management, 18(2), 1706-1717.

[11] Tursynova, A., & Omarov, B. (2021, November). 3D U-Net for brain stroke lesion segmentation on ISLES 2018 dataset. In 2021 16th International Conference on Electronics Computer and Computation (ICECCO) (pp. 1-4). IEEE.

[12] Asha, A., Arunachalam, R., Poonguzhali, I., Urooj, S., & Alelyani, S. (2023). Optimized RNN-based performance prediction of IoT and WSN-oriented smart city application using improved honey badger algorithm. Measurement, 210, 112505.

[13] Omarov, B., Altayeva, A., & Cho, Y. I. (2017). Smart building climate control considering indoor and outdoor parameters. In Computer Information Systems and Industrial Management: 16th IFIP TC8 International Conference, CISIM 2017, Bialystok, Poland, June 16-18, 2017, Proceedings 16 (pp. 412-422). Springer International Publishing.

[14] Rawal, B. S., Manogaran, G., Singh, R., Poongodi, M., & Hamdi, M. (2021, June). Network augmentation by dynamically splitting the switching function in SDN. In 2021 IEEE International Conference on Communications Workshops (ICC Workshops) (pp. 1-6). IEEE.

[15] Latif, S. A., Wen, F. B. X., Iwendi, C., Li-Li, F. W., Mohsin, S. M., Han, Z., & Band, S. S. (2022). AI-empowered, blockchain and SDN integrated security architecture for IoT network of cyber physical systems. Computer Communications, 181, 274-283.

[16] Wang, Y., Shang, F., Lei, J., Zhu, X., Qin, H., & Wen, J. (2023). Dual-attention assisted deep reinforcement learning algorithm for energy-

efficient resource allocation in Industrial Internet of Things. Future Generation Computer Systems, 142, 150-164.

[17] UmaMaheswaran, S. K., Prasad, G., Omarov, B., Abdul-Zahra, D. S., Vashistha, P., Pant, B., & Kaliyaperumal, K. (2022). Major challenges and future approaches in the employment of blockchain and machine learning techniques in the health and medicine. Security and Communication Networks, 2022.

[18] Cao, B., Sun, Z., Zhang, J., & Gu, Y. (2021). Resource allocation in 5G IoV architecture based on SDN and fog-cloud computing. IEEE Transactions on Intelligent Transportation Systems, 22(6), 3832-3840.

[19] Keshari, S. K., Kansal, V., Kumar, S., & Bansal, P. (2023). An intelligent energy efficient optimized approach to control the traffic flow in Software-Defined IoT networks. Sustainable Energy Technologies and Assessments, 55, 102952.

[20] Poornima, E., Muthu, B., Agrawal, R., Kumar, S. P., Dhingra, M., Asaad, R. R., & Jumani, A. K. (2023). Fog robotics-based intelligence transportation system using line-of-sight intelligent transportation. Multimedia Tools and Applications, 1-29.

[21] Razdan, S., & Sharma, S. (2022). Internet of medical things (IoMT): Overview, emerging technologies, and case studies. IETE technical review, 39(4), 775-788.

[22] Kazmi, S. H. A., Qamar, F., Hassan, R., Nisar, K., & Chowdhry, B. S. (2023). Survey on joint paradigm of 5G and SDN emerging mobile technologies: Architecture, security, challenges and research directions. Wireless Personal Communications, 1-48.

[23] Amiri, Z., Heidari, A., Navimipour, N. J., & Unal, M. (2023). Resilient and dependability management in distributed environments: A systematic and comprehensive literature review. Cluster Computing, 26(2), 1565-1600.

[24] Banafaa, M., Shayea, I., Din, J., Azmi, M. H., Alashbi, A., Daradkeh, Y. I., & Alhammadi, A. (2023). 6G mobile communication technology: Requirements, targets, applications, challenges, advantages, and opportunities. Alexandria Engineering Journal, 64, 245-274.

[25] Aboamer, M. A., Sikkandar, M. Y., Gupta, S., Vives, L., Joshi, K., Omarov, B., & Singh, S. K. (2022). An investigation in analyzing the food quality well-being for lung cancer using blockchain through cnn. Journal of Food Quality, 2022.

[26] Ray, P. P., & Kumar, N. (2021). SDN/NFV architectures for edge-cloud oriented IoT: A systematic review. Computer Communications, 169, 129-153.

[27] Naeem, F., Ali, M., & Kaddoum, G. (2023). Federated-learning-empowered semi-supervised active learning framework for intrusion detection in ZSM. IEEE Communications Magazine, 61(2), 88-94.

[28] Mughaid, A., AlZu'bi, S., Alnajjar, A., AbuElsoud, E., Salhi, S. E., Igried, B., & Abualigah, L. (2023). Improved dropping attacks detecting system in 5g networks using machine learning and deep learning approaches. Multimedia Tools and Applications, 82(9), 13973-13995.

[29] Rahman, A., Islam, M. J., Montieri, A., Nasir, M. K., Reza, M. M., Band, S. S., ... & Mosavi, A. (2021). Smartblock-sdn: An optimized blockchain-sdn framework for resource management in iot. IEEE Access, 9, 28361-28376.

[30] Ribeiro, D. A., Melgarejo, D. C., Saadi, M., Rosa, R. L., & Rodríguez, D. Z. (2023). A novel deep deterministic policy gradient model applied to intelligent transportation system security problems in 5G and 6G network scenarios. Physical Communication, 56, 101938.

[31] Javanmardi, S., Shojafar, M., Mohammadi, R., Persico, V., & Pescapè, A. (2023). S-FoS: A secure workflow scheduling approach for performance optimization in SDN-based IoT-Fog networks. Journal of Information Security and Applications, 72, 103404.

[32] Kashef, M., Visvizi, A., & Troisi, O. (2021). Smart city as a smart service system: Human-computer interaction and smart city surveillance systems. Computers in Human Behavior, 124, 106923.

[33] Qu, Y., Wang, Y., Ming, X., & Chu, X. (2023). Multi-stakeholder's sustainable requirement analysis for smart manufacturing systems based on the stakeholder value network approach. Computers & Industrial Engineering, 177, 109043.

[34] Bourechak, A., Zedadra, O., Kouahla, M. N., Guerrieri, A., Seridi, H., & Fortino, G. (2023). At the Confluence of Artificial Intelligence and Edge Computing in IoT-Based Applications: A Review and New Perspectives. Sensors, 23(3), 1639.

[35] Imam-Fulani, Y. O., Faruk, N., Sowande, O. A., Abdulkarim, A., Alozie, E., Usman, A. D., ... & Taura, L. S. (2023). 5G Frequency Standardization, Technologies, Channel Models, and Network Deployment: Advances, Challenges, and Future Directions. Sustainability, 15(6), 5173.

[36] Abou El Houda, Z., Hafid, A. S., & Khoukhi, L. (2023). Mitfed: A privacy preserving collaborative network attack mitigation framework based on federated learning using sdn and blockchain. IEEE Transactions on Network Science and Engineering.

[37] Sheng, M., Zhou, D., Bai, W., Liu, J., Li, H., Shi, Y., & Li, J. (2023). Coverage enhancement for 6G satellite-terrestrial integrated networks: performance metrics, constellation configuration and resource allocation. Science China Information Sciences, 66(3), 130303.

[38] Sutradhar, S., Karforma, S., Bose, R., & Roy, S. (2023). A Dynamic Step-wise Tiny Encryption Algorithm with Fruit Fly Optimization for Quality of Service improvement in healthcare. Healthcare Analytics, 3, 100177.

[39] Al-Turjman, F., Zahmatkesh, H., & Shahroze, R. (2022). An overview of security and privacy in smart cities' IoT communications. Transactions on Emerging Telecommunications Technologies, 33(3), e3677.

[40] Mahi, M. J. N., Chaki, S., Ahmed, S., Biswas, M., Kaiser, M. S., Islam, M. S., ... & Whaiduzzaman, M. (2022). A review on VANET research: Perspective of recent emerging technologies. IEEE Access, 10, 65760-65783.

[41] Ahmad, S., & Mir, A. H. (2021). Scalability, consistency, reliability and security in SDN controllers: a survey of diverse SDN controllers. Journal of Network and Systems Management, 29, 1-59.

[42] Zhou, H., Zheng, Y., Jia, X., & Shu, J. (2023). Collaborative prediction and detection of DDoS attacks in edge computing: A deep learning-based approach with distributed SDN. Computer Networks, 225, 109642.

[43] Zhang, J., Liu, Y., Li, Z., & Lu, Y. (2023). Forecast-assisted service function chain dynamic deployment for SDN/NFV-enabled cloud management systems. IEEE Systems Journal.

[44] Priyadarshini, R., & Barik, R. K. (2022). A deep learning based intelligent framework to mitigate DDoS attack in fog environment. Journal of King Saud University-Computer and Information Sciences, 34(3), 825-831.

[45] Das, S. K., Benkhelifa, F., Sun, Y., Abumarshoud, H., Abbasi, Q. H., Imran, M. A., & Mohjazi, L. (2023). Comprehensive review on ML-based RIS-enhanced IoT systems: basics, research progress and future challenges. Computer Networks, 224, 109581.

[46] Mubarakali, A., Durai, A. D., Alshehri, M., AlFarraj, O., Ramakrishnan, J., & Mavaluru, D. (2023). Fog-based delay-sensitive data transmission algorithm for data forwarding and storage in cloud environment for multimedia applications. Big Data, 11(2), 128-136.

[47] Liu, D., Li, Z., & Jia, D. (2023). Secure distributed data integrity auditing with high efficiency in 5G-enabled software-defined edge computing. Cyber Security and Applications, 1, 100004.

[48] Kazmi, S. H. A., Qamar, F., Hassan, R., & Nisar, K. (2023). Routing-based interference mitigation in SDN enabled beyond 5G communication networks: A comprehensive survey. IEEE Access.

[49] Gong, J., & Rezaeipanah, A. (2023). A fuzzy delay-bandwidth guaranteed routing algorithm for video conferencing services over SDN networks. Multimedia Tools and Applications, 1-30.

[50] Alani, T. O., & Al-Sadi, A. M. (2023). Survey of optimizing dynamic virtual local area network algorithm for software-defined wide area network. TELKOMNIKA (Telecommunication Computing Electronics and Control), 21(1), 77-87.

# Automatic Recognition of Marine Creatures using Deep Learning

Oudayrao Ittoo, Sameerchand Pudaruth

ICT Department, FoICDT, University of Mauritius, Réduit, Mauritius

*Abstract*—The identification of marine species is a challenge for people all over the world, and the situation is not different for Mauritians. It is of utmost importance to create an automated system to correctly identify marine species. In the past, researchers have used machine learning to address the issue of marine creature recognition. The manual feature extraction part of machine learning complicates model creation as features have to be extracted manually using an appropriate filter. In this work, we have used deep learning models to automate the feature extraction procedure. Currently, there is no publicly available dataset of marine creatures from the Indian Ocean. We created one of the biggest datasets used in this field, consisting of 51 different marine species collected from the Odysseo Oceanarium in Mauritius. The original dataset has a total of 5,709 images and is imbalanced. Image augmentations were performed to create an oversampled version of the dataset with 171 images per class, for a total of 8,721 images. The MobileNetV1 model trained on the oversampled dataset with a split ratio of 80% for training and 10% for validation and testing was the best performing one in terms of classification accuracy and inference time. The model had the smallest inference time of 0.10 seconds per image and attained a classification accuracy of 99.89% and an F1 score of 99.89%.

*Keywords—Marine creature identification; machine learning; deep learning; MobileNetV1; Mauritius*

## I. INTRODUCTION

Millions of marine creatures live in the ocean, making it the largest habitat on the planet [1]. The world's health is closely related to this marine biodiversity. These marine creatures are of the utmost importance to society as they are a source of food and a symbol of economic welfare. For example, fish are estimated to provide 20% of animal protein to about three billion people [2]. In addition, the ocean is home to a diverse range of creatures that can be utilised for the development of pharmaceutical products to treat various diseases [3]. The human race benefits from the numerous advantages that the marine ecosystem provides for its survival [4]. Therefore, the effective conservation of this biodiversity in a sustainable manner is crucial for the proper functioning of the marine ecosystem and the human race [5].

Traditionally, marine biologists identify aquatic species by visually inspecting their morphological traits [6-7]. Another popular method to correctly identify and group them is deoxyribonucleic acid (DNA) barcoding [8]. This method can be used to precisely, accurately, and quickly detect invasive alien species or marine bacteria that can cause viral outbreaks [9]. DNA barcoding has already proven its worth as a deterrent against various forms of economic fraud, such as seafood

mislabelling [10]. Despite its advantages, such an identification method is labour-intensive and time-consuming. As a result, it is crucial to create an automatic marine creature recognition system to address these difficulties.

Automatic recognition of marine creatures is a topic of interest to many researchers around the world. Pudaruth et al. have developed such a system in the form of a mobile application to recognise some of the marine fish that are present in Mauritian waters [11]. However, no system has been developed to cater for other types of marine species that can be found in the Indian Ocean. Common people with no expertise in marine taxonomy have difficulties distinguishing between the different aquatic species. This poses a problem, especially when deadly ocean animals such as the stonefish, blue-ringed octopus, or the lionfish, amongst many others, are encountered [12-13]. Furthermore, some endangered species require proper protection, such as conservation laws and regulations. These are only feasible after recognising them.

The motivation for this study is to create an image recognition model capable of distinguishing between different marine creatures. It is worth mentioning that 80% of the ocean is still undiscovered. According to an interview given by Dr Gene Carl Feldman to Oceana (an organisation focused on protecting the ocean), space exploration is far simpler than ocean exploration [14]. As there is an abundance of marine life in the ocean and it is difficult to cater for all of them, the scope of this study focuses only on some marine creatures that are available at Odysseo Oceanarium in Mauritius. The recognition model has been integrated into a web application. The importance of this application is diverse. First, it will help in creating awareness about the creatures, especially dangerous species. Furthermore, it will also help in raising knowledge by providing some basic information about the creature after the recognition phase. Information such as its scientific name, common name, short description, and whether the animal is deadly has been provided. The information provided can then be used to better understand the animal. The proposed system employs computer vision and deep learning techniques to properly and accurately identify the marine creatures.

This paper is divided into different sections. In Section II, a background study and reviews of related work in this field are provided. Section III delves deeper into the methodology and Section IV assesses the model's performance and discusses the results obtained. The final section concludes this paper. Table XII in Appendix I lists all the marine creatures from the Odysseo Oceanarium which were used in this study.

## II. LITERATURE REVIEW

Several studies have been carried out in the past to develop systems for the automated identification of marine life. This section contains summaries of several relevant works in chronological order.

### A. Fish Recognition

Strachan et al. conducted one of the earliest studies in this field, attempting to evaluate three different image analysis methods to differentiate between images of fish from different species [15]. Methods such as invariant moments, mismatch optimisation, and geometric shape descriptors were used. Their strategy takes into account the fact that fish can be identified by their body shape. The dataset used in their experiment consisted of 60 different fish images. Their research found that the geometric shape descriptors method outperformed the other two approaches, yielding a 90% accuracy rate. However, their experiment was limited to a restricted number of species: only seven different species (two of which were identical gurnard fish, shot from different perspectives).

Fish image recognition has found its usefulness in systems such as automated fish counting. Fish counting is a challenging but critical task for the maintenance of a sustainable fishing level and the prevention of overfishing [16]. Luo et al. proposed a method for such a system by using video footage captured during fishing operations [17]. Their method involved the use of Statistical Shape Models (SSM) and Artificial Neural Network (ANN). To overcome the occlusion problem caused by people walking around the deck of the ship, the video footage was pre-processed. The colour of the images was used as a feature for the recognition process, and an Error Back Propagating ANN classifier was used to recognise the fish from the background. The next step was to use SSM to identify the fish. Lastly, a rule-based counting method was used to count the number of fish. Their method achieved an accuracy of 89.6% for a one-hour video.

The study conducted by Rathi et al. proposed a solution based on Convolutional Neural Network (CNN), deep learning, and image processing [18]. Their method involved pre-processing of the captured images with the aim of removing noise and then using CNN to classify the different fishes. The Otsu's thresholding method was adopted to obtain a histogram representation of the input grayscale image. The next step was to perform morphological operations, such as dilation and erosion, to prepare the resulting image to be processed by the CNN algorithm. To put their method to the test, they used 27,142 images from the Fish4Knowledge dataset, representing 21 species which resulted in an accuracy of 96.29%. However, due to background noise and a lack of image enhancement techniques to compensate for lost features during the pre-processing phase, some of the classifications were incorrect.

Faster R-CNN was used by Mandal et al. to create a system for the automatic detection and identification of fish species [19]. Their dataset consisted of 50 different fish species. Using a random sample technique, their dataset was divided into training (70%), validation (10%), and testing (20%) sets. They were able to achieve a mean average precision (mAP) of 82.4%.

Deep and Dash employed CNN for feature extraction, followed by Support Vector Machine (SVM) and k-Nearest Neighbour (kNN) for classification [20]. They used the Fish4Knowledge dataset which was divided into training (90%) and testing (10%) sets. The training set was further divided such that 10% of the training images were used for validation. Their research proved that using kNN for classification yields the best accuracy of 98.79%.

Rico-Díaz et al. proposed a non-invasive method for addressing the fish recognition problem by combining artificial vision techniques and ANN [21]. Their work relied on the fact that fish from different species can be distinguished based on their eye's sclera and pupil. The first step in the identification process was to employ image filtering techniques to reduce noise in the captured image. After that, background subtraction was done to segment the fish from the background. The next step was to identify the fish's eye, and for this the Hough algorithm was used. Additionally, a feed-forward ANN, being more costly, was also employed if the first method (the Hough algorithm) failed. Using their approach, they were able to achieve an overall accuracy of 74% for eye detection. Also, two underwater cameras were used to estimate the size and weight of the fish while they were swimming. Their solution, however, is dependent on the image's quality and a good background subtraction. Furthermore, performance degrades when the ANN is used if the Hough algorithm fails to detect the fish.

Liang et al. combined CNN and migration learning to distinguish between three different kinds of Chinese ornamental fish [22]. They used TensorFlow, which is an open-source library for machine learning and artificial intelligence, to train their network model. A total of 14, 000 (4*3,500) images were gathered, which were divided into 3,000 and 500 images for training and testing, respectively. Their dataset, which consisted of 3,500 images of three different fish and one set for other types of fish, was gathered from the Internet by using the web crawler technology. Following that, pre-processing was an important step in enhancing the recognition rate as real-time videos of fish were shot outside of the aquarium. The dark channel prior and gamma correction methods were used for this purpose. The latter was a significant step towards the removal of brightness from the pictures. To reduce processing power, all images were scaled to 250 * 250 pixels. Their experiment showed that an accuracy of 98.1% is achievable.

Cai et al. took a different approach to realising a system for detecting fish and counting [23]. They proposed to use the You Only Look Once Version 3 (YOLOv3) model with MobileNet as the backbone for feature extraction. Their proposed system was trained using different strategies. They found out that their system performs better than when using YOLOv3 alone. The average precision obtained was 79.61%.

Pudaruth et al. experimented with multiple machine learning classifiers to discover the most effective one for developing a smartphone application to recognise different fish species existing in the exclusive economic zone (EEZ) of Mauritius [11]. Their model was tested on 38 different fish species with a dataset consisting of 1,520 images. Using the

kNN classifier, they were able to attain an accuracy of 96%. Also, their study has shown that the use of deep neural networks (DNN) with the TensorFlow framework can attain an impressive accuracy of 98%. However, pictures of the fish were taken in a controlled environment and not in their natural habitat. The fish were placed on a white background to enable easier segmentation.

Conrady et al. used the Mask Region-Based Convolutional Neural Network (R-CNN) object detection framework to perform classification of the Roman seabream fish, which is endemic in Southern Africa [24]. Their dataset consisted of 2,015 images of the fish. They were able to get a mAP of 81.45% on their test data.

### B. Marine Creature Recognition

While most researchers have focused on fish recognition, Chen and Yu proposed a high-definition camera system capable of recognising marine creatures [25]. Their approach to the identification process was broken down into multiple steps. The first step was concerned with the extraction of the image frames from the original video. The following step involved pre-processing of the retrieved images. In addition, for the detection phase, the original image had to be transformed to its grayscale and binary representations. To classify seven creatures, two separate methods were used: The Back Propagation Neural Network and the SVM methods. Their study showed that the SVM approach had an accuracy of 92% in classifying the creatures, which was higher compared to the other methods. However, their proposed method does not work well for creatures with similar shapes.

Pelletier et al. developed a system capable of classifying marine animals into eight categories [26]. Their imbalanced dataset contained 3,777 images. They used two models to conduct their tests, namely AlexNet and GoogLeNet. The best performing model was GoogLeNet. The models were tested with uncropped and cropped images. They found out that by using the cropped images and GoogLeNet, they got the best accuracy, which was 96.54%. Additional tests were done by also considering the top two results during the classification process. This further increased the accuracy of the GoogLeNet model with the cropped dataset to 98.94%. Aside from image recognition, several other approaches for automatic identification of marine creatures have been used in the recent past. Demertzis et al. suggested a novel technique: the use of a Machine Hearing Framework (MHF) for the identification of marine animals through their underwater sounds [27]. They were able to recognise fish and marine animals with recognition accuracy of 96.08% and 92.18%, respectively.

Song et al. proposed a method for the identification of marine creatures from seafloor videos [28]. Their proposed methods are twofold: extraction of valid video clips followed by recognition of the creature. During the first phase, an image segmentation method was used to determine and extract all frames from the video containing the marine creature. The next phase was concerned with identifying and labelling the creatures in the valid video clips. This was accomplished with the help of public participation. Lastly was the recognition process, which was accomplished with the help of the information submitted by the public and the membership function. Their method had an accuracy of above 80% in extracting the valid video frames, and all the creatures were successfully recognised.

Liu et al. implemented an embedded system to classify marine animals into seven categories [29]. Their dataset includes 8,455 photos of marine animals, with 80% of the images used for training and 20% used for validation. The training images were augmented by applying some transformations (rotation, translation, and flipping), which increased the number of training images to 27,056 (6,764*4). The models that were deployed on the embedded device were tested with 350 new images. Three models were used: MobileNetV1, MobileNetV2, and InceptionV3. MobileNetV2 had the highest testing accuracy of 95.0% and validation accuracy of 92.89%. Their model took an average of 0.0578 seconds to classify one image.

### C. Knowledge Gap

Even though multiple studies have been conducted, very few have tested the effect of using deep learning (DL) on a big dataset. According to the review, the largest dataset had 50 species and consisted of 4,909 images [19]. Furthermore, there is no web application available in Mauritius that can perform recognition of marine animals. Adding to that, no dataset consisting of more than 50 marine creature species is currently available. This work aims to provide some answers to these research gaps as well as to contribute a dataset and a web application to perform marine creature classification.

## III. METHODOLOGY

### A. Data Collection at Odysseo Oceanarium

As no relevant existing dataset of marine creatures from the Indian Ocean was found at the time of this study, a custom dataset was created. Data collection can be done from multiple sources. However, due to time constraints, this is not feasible. Nonetheless, the source should be trustworthy and provide the desired information. In this regard, the Odysseo Oceanarium was chosen as the primary source of data gathering for this study. In recent years, underwater video surveillance has grown increasingly common in marine environments to acquire data on marine creatures in their natural habitat. This is a non-invasive method and provides sufficient data for research. Chen and Yu adopted this approach by using an underwater submerged video system [25]. However, for this study, due to limited resources, videos of marine creatures were taken outside of the aquariums found at Odysseo using a smartphone. All the videos were taken with a Huawei Y9 Prime 2019 smartphone, which has a resolution of 16 megapixels.

### B. Data Processing

Numerous videos of marine creatures were obtained at Odysseo. Each video was converted into frames. Pelletier et al. have already shown that cropped images result in better model performance [26]. Taking this into consideration, each extracted frame was carefully cropped. Not all images were included in the final dataset. The following conditions had to be met to use the image, or else it was discarded: the image containing the creature should not be occluded by another creature or object; the creature should be recognisable and it

should not be too far from the image. Fig. 1 shows an example of a good image. It follows all the criteria described above.



Fig. 1.    Example of a good image.

## C. Custom Dataset

The custom-built dataset consists of 51 classes of marine creatures as shown in Fig. 2. It has 5,709 images in total and is imbalanced. The class having the highest number of images is Dascyllus aruanus with 171 pictures, and the one with the lowest number of images is Chaetodon kleinii with 74 pictures.



Fig. 2.    Unbalanced dataset distribution.

Marine animals are challenging to photograph since they conceal their presence in aquarium by hiding beneath rocks, among plants and tank accessories. This makes it difficult to collect the same number of images for all of the creatures and is the primary reason why the dataset is initially imbalanced.

## D. Oversampling

To achieve an equal distribution of images per class, the dataset had to be balanced and for this oversampling was done as shown in Fig. 3. The following transformations were used for augmentation: flipping, change in brightness, shearing and rotation. Each image is subject to three possible modifications.

After oversampling, all classes got an equal distribution of images. Each marine creature now has 171 images. Fig. 4 shows the data distribution of the dataset after oversampling.



Fig. 3.    Distribution after oversampling.

## E. Splitting the Dataset

The oversampled dataset (171 images per class) were split using splitting ratios of 8:1:1, 7:2:1 and 6:3:1 for training, validation and testing sets. As a result of multiple manipulations, different dataset versions were created. A proper naming convention was devised to properly organise the work being done. Table I lists the various datasets that were used throughout this paper.

TABLE I.        SUMMARY OF DATASET VERSIONS

| # | Dataset Name | Description |
|---|--------------|-------------|
| 1 | DS_oversample_8_1_1 | This is the oversampled dataset, which contains 171 images for each class. Three different splits are done. |
| 2 | DS_oversample_7_2_1 | |
| 3 | DS_oversample_6_3_1 | |

## F. Feature Extraction and Classification

For this research, pre-trained CNN models were used both for feature extraction and classification. The pre-trained models that were employed are: MobileNetV1, InceptionV3 and VGG16. Different image sizes were utilised depending on which model was imported. The image sizes used are shown in Table II [30-31].

TABLE II.        IMAGE SIZES FOR DIFFERENT MODELS

| Model Name | Image Size |
|------------|------------|
| MobileNetV1 | 224*224 |
| InceptionV3 | 299*299 |
| VGG16 | 224*224 |

Fig. 4. Oversampling flowchart.

After the features have been extracted, the next step is to use these features and a classifier to make a prediction. The architecture of the pre-trained models used is shown in Fig. 5. The input image given to the VGG16 and MobileNetV1 models is an image of size 224 * 224 compared to the InceptionV3 model, which is of size 299 * 299.

The classification layers of the pre-trained model were replaced with one global average pooling layer and three dense layers. The softmax activation function was used in the model's final dense layer for classification. Adding to that, for the two other dense layers, the rectified linear activation unit (ReLU) activation function was used. The custom model predicts the input image as one of the 51 classes of marine creatures from the dataset.

Fig. 5.   High level architecture of pre-trained model.

### G.  Use of Callback Functions

Different callback functions, such as ModelCheckpoint and EarlyStopping, were employed during model training.  The EarlyStopping callback function is viewed as a technique to combat model overfitting. For this work, the number of epochs was fixed at 100, and then the EarlyStopping function was used to halt training if the model became overfit [32]. The ModelCheckpoint callback function, on the other hand, was used to save the model during the training phase. Failure may occur occasionally, causing the training to be disrupted. It is preferable to resume training from the last saved epoch rather than starting it from scratch [33].

### H.  Use of Optimiser

The Stochastic Gradient Descent (SGD) optimiser was employed during model training to adjust the weight and learning rate properties of the DL model to reduce losses during backpropagation. Additionally, to help the optimiser converge in the right direction and prevent overshooting, Nesterov momentum was employed. Due to its look ahead property, the Nesterov method takes the appropriate precautions by making smaller updates to reach the minima [34-35]. Experiments were performed using the MobileNetV1 pre-trained model to find the optimal parameters to pass to the SGD optimiser function. We discovered that setting the learning rate to 0.001, decay to 1e-6, and momentum to 0.8 produces the best results. Thus, these parameter values were used throughout this work for model training.

### I.  Training using Transfer Learning

Transfer learning is a concept whereby a previously trained (pre-trained) model is reused to tackle a new but comparable task. This is a popular deep learning technique as the neural network does not have to be trained from scratch with a huge volume of data. The weight that the network has already learnt is simply transferred to the new task in transfer learning. This technique aids in reducing training time and may possibly increase the performance of the neural network [36]. In this research, transfer learning is used to train the CNN models. Fig. 6 illustrates how transfer learning was applied for model training using the custom-built dataset.



Fig. 6.   Block diagram of the proposed training strategy.

There are different types of fine-tuning that can be done on pre-trained CNN models, such as training the entire model, training some layers and leaving others frozen, and freezing the convolutional base by not training the feature extraction layers [37]. For this work, each test will be done by training the entire model and freezing the convolutional base.

### J.  The Web Application

A good and simple user interface is crucial for the user to efficiently use the application. Taking this into consideration, the user interface of the web application is divided into 4 areas: image upload area, prediction area, creatures in dataset area and modal displaying creature information.

*1)  Image upload area:* Fig. 7 shows the image upload area when the application is accessed through a desktop and a mobile phone.



Fig. 7.   Image upload area for desktop view (left) and mobile view (right).

*2) Prediction area:* The prediction area is illustrated in Fig. 8. It shows the predicted creature, along with a table containing creatures with confidence scores greater than the threshold value.



Fig. 8. Prediction area (system can categorize the image).

If the confidence scores computed for the input image are lower than the threshold, an appropriate message is displayed to the user as shown in Fig. 9. In this context, a threshold value of 0.5 is employed.



Fig. 9. Prediction area (system cannot categorize the image).

*3) Creatures in dataset area:* The creatures that the system can classify are shown in Fig. 10.



Fig. 10. DataTable for displaying creatures from the dataset.

*4) Modal displaying creature information:* The modal component is used to display information about a creature as shown in Fig. 11.



Fig. 11. Modal displaying creature information.

## IV. RESULTS

This section provides an evaluation of the different models.

### A. Testing Different Split Ratios

Table III shows the three splits for the oversampled dataset.

TABLE III. OVERSAMPLED DATASET SPLIT SUMMARY

| Dataset | Split Ratio | Train | Val | Test |
|---|---|---|---|---|
| DS_oversample_8_1_1 | Train: 80% Val: 10% Test: 10% | 136 | 17 | 18 |
| DS_oversample_7_2_1 | Train: 70% Val: 20% Test: 10% | 119 | 34 | 18 |
| DS_oversample_6_3_1 | Train: 60% Val: 30% Test: 10% | 102 | 51 | 18 |

The three pre-trained models employed were trained on the three versions of the oversampled dataset. Table IV, Table V and Table VI shows the results obtained.

TABLE IV. DS_OVERSAMPLE_8_1_1 RESULT

| Model | Trainable Layers | Classification Accuracy (%) | F1 Score (%) |
|---|---|---|---|
| MobileNetV1 | False | 99.56 | 99.56 |
| | True | 99.89 | 99.89 |
| VGG16 | False | 98.58 | 98.58 |
| | True | 98.91 | 98.91 |
| InceptionV3 | False | 97.17 | 97.16 |
| | True | 99.35 | 99.35 |

TABLE V. DS_Oversample_7_2_1 Result

| Model | Trainable Layers | Classification Accuracy (%) | F1 Score (%) |
|---|---|---|---|
| MobileNetV1 | False | 99.67 | 99.68 |
| | True | 99.89 | 99.89 |
| VGG16 | False | 98.47 | 98.48 |
| | True | 99.89 | 99.89 |
| InceptionV3 | False | 97.60 | 97.58 |
| | True | 99.46 | 99.45 |

TABLE VI. DS_Oversample_6_3_1 Result

| Model | Trainable Layers | Classification Accuracy (%) | F1 Score (%) |
|---|---|---|---|
| MobileNetV1 | False | 98.58 | 98.58 |
| | True | 99.46 | 99.45 |
| VGG16 | False | 97.82 | 97.82 |
| | True | 99.02 | 99.03 |
| InceptionV3 | False | 96.51 | 96.52 |
| | True | 99.35 | 99.35 |

TABLE VII. Best Models for Oversampled Dataset

| Dataset | Split Ratio | Best Model | Classification Accuracy (%) | F1 Score (%) |
|---|---|---|---|---|
| DS_oversample_8_1_1 | Train: 80% Val: 10% Test: 10% | MobileNetV1 | 99.89 | 99.89 |
| DS_oversample_7_2_1 | Train: 70% Val: 20% Test: 10% | MobileNetV1 | 99.89 | 99.89 |
| | | VGG16 | 99.89 | 99.89 |
| DS_oversample_6_3_1 | Train: 60% Val: 30% Test: 10% | MobileNetV1 | 99.46 | 99.45 |

The best accuracy obtained in all cases is when the feature extraction layers are trained. Table VII shows the best performing model for each of the different splits.

Irrespective of the split ratios used, the best models were able to achieve very good accuracy of above 99%. The variations in the scores, as indicated in Table VII are less than 1%. This is due to the fact that randomness is used in weight initialization when training of the model starts. The weights are adjusted at every epoch. This produces different outcomes for the same model each time it is trained on the same dataset [38-39]. This means that if the same experiments were repeated, different scores would have been obtained. To conclude, the differences obtained are considered insignificant.

Judging the models solely on accuracy is not enough to give a fair evaluation. The model's prediction time must also be considered. The inference time of the models to predict all the images in their test directory was repeated five times. In Table VIII, the total time taken is displayed in seconds.

Using the values shown in Table VIII the average inference time can be calculated by dividing the time taken to predict all the images by five. The prediction time for one image can then

be computed by dividing the resulting value by the number of test pictures as shown in Table IX.

TABLE VIII. Prediction Time

| Dataset | Number of Images | MobileNetV1 | VGG16 | InceptionV3 |
|---|---|---|---|---|
| DS_oversample_8_1_1 | 918 (18*51) | 475.5 | 2453.5 | 1002.6 |

TABLE IX. Prediction Time Per Image

| Dataset | MobileNetV1 | VGG16 | InceptionV3 |
|---|---|---|---|
| DS_oversample_8_1_1 | 0.10 | 0.53 | 0.22 |

From Table IX, it can be seen that the MobileNetV1 models achieved the lowest inference time. Table X shows the number of trainable parameters and the sizes of the three models.

TABLE X. Model Sizes

| Model | Parameters | Size |
|---|---|---|
| VGG16 | Total params: 14,885,491<br>Trainable params: 14,885,491<br>Non-trainable params: 0 | 113.7 MB |
| InceptionV3 | Total params: 22,366,803<br>Trainable params: 22,332,371<br>Non-trainable params: 34,432 | 171.8 MB |
| MobileNetV1 | Total params: 3,530,739<br>Trainable params: 3,508,851<br>Non-trainable params: 21,888 | 27.3 MB |

Among the three models, MobileNetV1 has the least number of trainable parameters. Furthermore, the MobileNetV1 model is smaller in terms of size. As a result, it has the shortest inference time. For deployment, the MobileNetV1 model trained on the DS_oversample_8_1_1 dataset was chosen since it had the lowest inference time of 0.10 seconds per image.

### B. Comparisons with Related Works

Even though the accuracies obtained in this study cannot be truly compared with other researchers because the same datasets were not employed, an attempt to compare our work with previous studies is made in this section.

Deep and Dash conducted several experiments using a dataset of 23 creatures [20]. They used CNN for both feature extraction and classification. Additionally, they used a hybrid strategy in which CNN was used to extract features and a classifier (kNN or SVM) was used to classify them. They got the best accuracy of 98.79% when they used their custom-made CNN for feature extraction and kNN as a classifier. However, in this study, we were able to achieve higher accuracy when the pre-trained CNN models were used for both extraction and classification.

The training methodology adopted by Liu et al. is the same as the one we have used [29]. They also used transfer learning to train their models. They used the MobileNetV1,

MobileNetV2 and InceptionV3 models to perform feature extraction and classification. Table XI shows a comparison between the best performing model in their work and in ours.

TABLE XI.    COMPARISION OF THE BEST PERFORMING MODEL

| Work | Best Model | Classification Accuracy (%) | Inference Time (Seconds) |
|---|---|---|---|
| Ittoo and Pudaruth | MobileNetV1 trained on the DS_oversample_8_1_1 dataset | 99.89 | 0.10 |
| Liu et al. (2019) | MobileNetV2 | 95.0 | 0.0578 |

The MobileNetV2 model presented by Liu et al. is limited to predicting 7 species [29]. However, the model presented in this study can perform classification between 51 creatures. The larger the number of classes in a deep learning model, the more time the model generally takes to predict the image.

## V. CONCLUSIONS

Automatic recognition of marine creatures is a topic of interest to many researchers around the world. People who are unfamiliar with marine taxonomy have trouble discriminating between different aquatic organisms. This poses a problem, especially when deadly ocean animals are encountered. Several studies have been undertaken over the last few decades, but not many have examined the effect of training deep learning pre-trained CNN models on a large dataset.

In this research, three deep learning models, namely, MobileNetV1, InceptionV3, and VGG16, were investigated and implemented for the task of marine creature classification. To achieve the objectives of this study, a customised dataset of 51 available creatures from the Indian Ocean was built and used for training and testing the effectiveness of models. Images of marine creatures were collected at Odysseo Oceanarium in Mauritius.

Several experiments with different split ratios were carried out. The splits for the training and validation sets were varied, and that for the testing set was fixed at 10%. Transfer learning was used, and the models were fine-tuned by replacing their classification layers with new ones. Adding to that, two experiments were performed for each model: training the feature extraction layers and not training them. All of these tests were carried out in order to determine the optimal split ratio, dataset, and model.

It has been concluded that the best suited model was MobileNetV1 trained with the oversampled dataset with a split ratio 80% for training, 10% for validation and 10% for testing. The model attained a classification accuracy and an F1 score of 99.89%. The model had an inference time of 0.10 seconds per image. This model was then integrated into a web application.

Our research has thus demonstrated that deep learning models offer enormous potential for automating the process of marine creature recognition. The developed web application with the integrated MobileNetV1 model provides a reliable and fully automated tool for the classification of marine creatures without the need for expert assistance.

## REFERENCES

[1] BBC Earth. "Blue Planet II: The Prequel" [Video]. YouTube, 27 Sep. 2017. https://www.youtube.com/watch?v=_38JDGnr0vA.

[2] UN Food and Agriculture Organization. "7 reasons why we need to act now to #SaveOurOcean." Medium, 5 Jun. 2017. https://medium.com/@UNFAO/7-reasons-why-we-need-to-act-now-to-saveourocean-70671e3da38d.

[3] Newman, T. "Why scientists are searching the ocean for new drugs." Medical News Today, 8 Jun. 2019. https://www.medicalnewstoday.com/articles/325384#Why-look-to-the-sea?

[4] Stuchtey, M. R., Vincent, A., Merkl, A., Bucher, M., Haugan, P. M., Lubchenco, J., & Pangestu, M. E. "Ocean Solutions That Benefit People, Nature and the Economy." World Resources Institute, 2020. https://oceanpanel.org/wp-content/uploads/2022/06/full-report-ocean-solutions-eng.pdf

[5] Luypaert, T., Hagan, J. G., McCarthy, M. L., & Poti, M. "Status of Marine Biodiversity in the Anthropocene." In S. Jungblut, V. Liebich & M. Bode-Dalby (Eds.), YOUMARES 9 - The Oceans: Our Research, Our Future, pp. 57–82. Springer, Cham, 2020. https://doi.org/10.1007/978-3-030-20389-4_4.

[6] Buckley, M., Fraser, S., Herman, J., Melton, N. D., Mulville, J., & Pálsdóttir, A. H. "Species identification of archaeological marine mammals using collagen fingerprinting." Journal of Archaeological Science, vol. 41, 2014, pp. 631-641. https://doi.org/10.1016/j.jas.2013.08.021.

[7] Goodwin, K. D., Thompson, L. R., Duarte, B., Kahlke, T., Thompson, A. R., Marques, J. C., & Caçador, I. "DNA Sequencing as a Tool to Monitor Marine Ecological Status." Frontiers in Marine Science, vol. 4, 2017, p. 107. https://doi.org/10.3389/fmars.2017.00107.

[8] Trivedi, S., Aloufi, A. A., Ansari, A. A., & Ghosh, S. K. "Role of DNA barcoding in marine biodiversity assessment and conservation: An update." Saudi Journal of Biological Sciences, vol. 23, no. 2, 2020, pp. 167-171. http://dx.doi.org/10.1016/j.sjbs.2015.01.001.

[9] Thomas, V. G., Hanner, R. H., & Borisenko, A. V. "DNA-based identification of invasive alien species in relation to Canadian federal policy and law, and the basis of rapid-response management." Genome, vol. 59, no. 11, 2016, pp. 1023-1031. https://doi.org/10.1139/gen-2016-0022.

[10] Harris, D. J., Rosado, D., & Xavier, R. "DNA Barcoding Reveals Extensive Mislabeling in Seafood Sold in Portuguese Supermarkets." Journal of Aquatic Food Product Technology, vol. 25, no. 8, 2016, pp. 1375-1380. http://dx.doi.org/10.1080/10498850.2015.1067267.

[11] Pudaruth, S., Nazurally, N., Appadoo, C., Kishnah, S., & Chady, F. "SuperFish: A Mobile Application for Fish Species Recognition Using Image Processing Techniques and Deep Learning." International Journal of Computing and Digital Systems, vol. 10, no. 1, 2021, pp. 1157-1165. http://dx.doi.org/10.12785/ijcds/1001104.

[12] Proctor, K. "Burneside man stung by 'deadly' fish on honeymoon." The Westmorland Gazette, 20 Oct. 2011. https://www.thewestmorland gazette.co.uk/news/9314711.burneside-man-stung-by-deadly-fish-on-honeymoon/.

[13] Shersby, M. "13 of the most venomous sea creatures lurking in the water." Live Science, 28 Aug. 2023. https://www.livescience.com/animals/13-of-the-most-venomous-sea-creatures-on-earth/

[14] Petsko, E. "Why does so much of the ocean remain unexplored and unprotected?" Oceana, 8 Jun. 2020. https://oceana.org/blog/why-does-so-much-ocean-remain-unexplored-and-unprotected/

[15] Strachan, N. J. C., Nesvadba, P., & Allen, A. R. "Fish Species Recognition by Shape Analysis of Images." Pattern Recognition, vol. 23, no. 5, 1990, pp. 539-544. https://doi.org/10.1016/0031-3203(90)90074-U

[16] Chepkemoi, J. "How Many Fish Live In The Ocean?" WorldAtlas, 25 Apr. 2017. https://www.worldatlas.com/articles/how-many-fish-are-there-in-the-ocean.html

[17] Luo, S., Li, X., Wang, D., Li, J., & Sun, C. "Automatic Fish Recognition and Counting in Video Footage of Fishery Operations." 2015 International Conference on Computational Intelligence and Communication Networks (CICN), 2015, pp. 296-299. https://doi.org/10.1109/CICN.2015.66

[18] Rathi, D., Jain, S., & Indu, S. "Underwater Fish Species Classification using Convolutional Neural Network and Deep Learning." 2017 Ninth International Conference on Advances in Pattern Recognition (ICAPR), 2017, pp. 1-6. https://doi.org/10.1109/ICAPR.2017.8593044

[19] Mandal, R., Connolly, R. M., Schlacher, T. A., & Stantic, B. "Assessing fish abundance from underwater video using deep neural networks." 2018 International Joint Conference on Neural Networks (IJCNN), 2018, pp. 1-6. https://doi.org/10.1109/IJCNN.2018.8489482

[20] Deep, B. V., & Dash, R. "Underwater Fish Species Recognition Using Deep Learning Techniques." 2019 6th International Conference on Signal Processing and Integrated Networks (SPIN), 2019, pp. 665-669. https://doi.org/10.1109/SPIN.2019.8711657

[21] Rico-Díaz, Á. J., Rabuñal, J. R., Gestal, M., Mures, O. A., & Puertas, J. "An Application of Fish Detection Based on Eye Search with Artificial Vision and Artificial Neural Networks." Water, vol. 12, no. 11, 2020, p. 3013. https://doi.org/10.3390/w12113013

[22] Liang, J., Fu, Z., Lei, X., Dai, X., & Lv, B. "Recognition and Classification of Ornamental Fish Image Based on Machine Vision." 2020 International Conference on Intelligent Transportation, Big Data & Smart City (ICITBS), 2020, pp. 910-913. https://doi.org/10.1109/ICITBS49701.2020.00201

[23] Cai, K., Miao, X., Wang, W., Pang, H., Liu, Y., & Song, J. "A modified YOLOv3 model for fish detection based on MobileNetv1 as backbone." Aquacultural Engineering, vol. 91, 2020, p. 102117. https://doi.org/10.1016/j.aquaeng.2020.102117

[24] Conrady, C. R., Er, Ş., Attwood, C. G., Roberson, L. A., & Vos, L. d. "Automated detection and classification of southern African Roman seabream using mask R-CNN." Ecological Informatics, vol. 69, 2022, p. 101593. https://doi.org/10.1016/j.ecoinf.2022.101593

[25] Chen, Q., & Yu, H. "Deep sea high-definition camera system based on marine creature classification technology." 2015 IEEE 16th International Conference on Communication Technology (ICCT), 2015, pp. 78-81. https://doi.org/10.1109/ICCT.2015.7399797

[26] Pelletier, S., Montacir, A., Zakari, H., & Akhloufi, M. "Deep Learning for Marine Resources Classification in Non-Structured Scenarios: Training vs. Transfer Learning." 2018 IEEE Canadian Conference on Electrical & Computer Engineering (CCECE), 2018, pp. 1-4. https://doi.org/10.1109/CCECE.2018.8447682

[27] Demertzis, K., Iliadis, L. S., & Anezakis, V. D. "Extreme deep learning in biosecurity: the case of machine hearing for marine species identification." Journal of Information and Telecommunication, vol. 2, no. 4, 2018, pp. 492-510. https://doi.org/10.1080/24751839.2018.1501542

[28] Song, X., Guo, Y., Zhang, F., Wang, S., Yang, J., & Chang, Y. "Research on marine creature recognition from seafloor videos." OCEANS 2019 - Marseille, France, 2019, pp. 1-4. https://doi.org/10.1109/OCEANSE.2019.8867479

[29] Liu, X., Jia, Z., Hou, X., Fu, M., Ma, L., & Sun, Q. "Real-time Marine Animal Images Classification by Embedded System Based on Mobilenet and Transfer Learning." OCEANS 2019 - Marseille, France, 2019, pp. 1-5. https://doi.org/10.1109/OCEANSE.2019.8867190

[30] Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M. & Adam, H. "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications." arXiv, 2017. https://doi.org/10.48550/arXiv.1704.04861

[31] Rosebrock, A. "ImageNet: VGGNet, ResNet, Inception, and Xception with Keras." PyImageSearch, 20 Mar. 2017. https://pyimagesearch.com/2017/03/20/imagenet-vggnet-resnet-inception-xception-keras/

[32] Chen, B. "Early Stopping in Practice: an example with Keras and TensorFlow 2.0." Towards Data Science, 29 Jul. 2020. https://towardsdatascience.com/a-practical-introduction-to-early-stopping-in-machine-learning-550ac88bc8fd

[33] Janapati, V. "A High Level Overview of Keras ModelCheckpoint Callback." Medium, 31 Oct. 2020. https://medium.com/swlh/a-high-level-overview-of-keras-modelcheckpoint-callback-deae8099d786

[34] Chandra, A. L. "Learning Parameters, Part 2: Momentum-Based & Nesterov Accelerated Gradient Descent." Towards Data Science, 15 May 2019. https://towardsdatascience.com/learning-parameters-part-2-a190bef2d12

[35] Doshi, S. "Various Optimization Algorithms For Training Neural Network." Towards AI, 13 Jan. 2019. https://towardsdatascience.com/optimizers-for-training-neural-network-59450d71caf6

[36] Sharma, P. "Understanding Transfer Learning for Deep Learning." Analytics Vidhya, 1 Mar. 2021. https://www.analyticsvidhya.com/blog/2021/10/understanding-transfer-learning-for-deep-learning/

[37] Marcelino, P. "Transfer learning from pre-trained models." Towards Data Science, 23 Oct. 2018. https://towardsdatascience.com/transfer-learning-from-pre-trained-models-f2393f124751

[38] Brownlee, J. "How to Get Reproducible Results with Keras." Machine Learning Mastery, 14 Jun. 2017. https://machinelearningmastery.com/reproducible-results-neural-networks-keras/

[39] Sharma, P. "Basic Introduction to Convolutional Neural Network in Deep Learning." Analytics Vidhya, 1 Mar. 2022. https://www.analyticsvidhya.com/blog/2022/03/basic-introduction-to-convolutional-neural-network-in-deep-learning/

APPENDIX I

TABLE XII. LIST OF ALL MARINE CREATURES

| # | Common Name | Scientific Name | Image | Number of Images |
|---|---|---|---|---|
| 1 | Scissortail sergeant | Abudefduf sexfasciatus | | 115 |
| 2 | Doubleband surgeonfish | Acanthurus tennentii | | 112 |

| 3 | Convict surgeonfish | Acanthurus triostegus |  | 116 |
| 4 | Yellowfin surgeonfish | Acanthurus xanthopterus |  | 122 |
| 5 | Yellowbreasted wrasse | Anampses twistii |  | 84 |
| 6 | White-spotted puffer | Arothron hispidus |  | 141 |
| 7 | Guineafowl puffer | Arothron meleagris |  | 101 |
| 8 | Blackspotted puffer | Arothron nigropunctatus |  | 128 |
| 9 | Common jellyfish | Aurelia aurita |  | 100 |

| 10 | Whitespotted filefish | Cantherhines dumerilii |  | 132 |
| 11 | Valentin's sharpnose puffer | Canthigaster valentini |  | 100 |
| 12 | Blacktip reef shark | Carcharhinus melanopterus |  | 106 |
| 13 | Leopard hind | Cephalopholis leopardus |  | 146 |
| 14 | Sunburst butterflyfish | Chaetodon kleinii |  | 74 |
| 15 | Raccoon butterflyfish | Chaetodon lunula |  | 92 |

| 16 | Cushion star | Culcita sp |  | 111 |
|----|-------------|-----------|------|-----|
| 17 | Whitetail dascyllus | Dascyllus aruanus |  | 171 |
| 18 | Thorntail stingray | Dasyatis thetidis |  | 111 |
| 19 | Red starfish | Echinaster sepositus |  | 93 |
| 20 | Common mushroom coral | Fungia fungites |  | 111 |
| 21 | Undulated moray | Gymnothorax undulatus |  | 136 |
| 22 | Pennant coralfish | Heniochus acuminatus |  | 85 |

| 23 | Spotted seahorse | Hippocampus kuda |  | 76 |
|----|------------------|------------------|--------------------|-----|
| 24 | Black sea cucumber | Holothuria atra |  | 83 |
| 25 | Longhorn cowfish | Lactoria cornuta |  | 154 |
| 26 | Bengal snapper | Lutjanus bengalensis |  | 100 |
| 27 | Blacktail snapper | Lutjanus fulvus |  | 114 |
| 28 | Common bluestripe snapper | Lutjanus kasmira |  | 114 |
| 29 | Pinecone soldierfish | Myripristis murdjan |  | 127 |

| 30 | Big blue octopus | Octopus cyanea |  | 111 |
| 31 | Peacock mantis shrimp | Odontodactylus scyllarus |  | 113 |
| 32 | Red-toothed triggerfish | Odonus niger |  | 111 |
| 33 | Yellow boxfish | Ostracion cubicum |  | 170 |
| 34 | Painted spiny lobster | Panulirus versicolor |  | 133 |
| 35 | Orangetail filefish | Pervagor aspricaudus |  | 82 |

| 36 | Longfin batfish | Platax teira |  | 129 |
|----|----|----|----|----|
| 37 | Emperor angelfish | Pomacanthus imperator |  | 113 |
| 38 | Devil firefish | Pterois miles |  | 112 |
| 39 | White-banded triggerfish | Rhinecanthus aculeatus |  | 88 |
| 40 | Wedge-tail triggerfish | Rhinecanthus rectangulus |  | 95 |
| 41 | Giant guitarfish | Rhynchobatus djiddensis |  | 112 |

| 42 | False stonefish | Scorpaenopsis diabolus |  | 79 |
| 43 | Zebra shark | Stegostoma fasciatum |  | 129 |
| 44 | Greenfish sea cucumber | Stichopus chloronotus |  | 79 |
| 45 | Small spotted dart | Trachinotus bailloni |  | 100 |
| 46 | Fluted giant clam | Tridacna squamosa |  | 148 |
| 47 | Striped sea urchin | Tripneustes gratilla |  | 106 |
| 48 | Blueband goby | Valenciennea strigata |  | 112 |

| 49 | Moorish idol | Zanclus cornutus |  | 118 |
| 50 | Indian sail-fin surgeonfish | Zebrasoma desjardinii (Juvenile) |  | 112 |
| 51 | Twotone tang | Zebrasoma scopas |  | 102 |

# Enhanced Linear Regression Models for Resource Usage Prediction in Dynamic Cloud Environments

Xiaoxiao Ma

School of Transportation, Chongqing Vocational College of Transportation, Chongqing 402247, China

*Abstract*—In response to the diverse resource utilization patterns observed across enterprises, this study proposes the utilization of adaptable cloud services. A novel system framework is presented, capturing and logging resource consumption at discrete intervals. Subsequently, this recorded data serves as input for a linear regression model, functioning as a machine learning tool to predict resource utilization in forthcoming intervals, leveraging historical data stored within the regression module. To bolster the resilience of the linear regression model, various effective meta-heuristic techniques are integrated alongside the conventional linear regression methodology, facilitating more accurate anticipation of overloaded or under-loaded resource conditions before their occurrence in real-world scenarios. Simulations demonstrate that the hybrid algorithm, named Whale Optimization Algorithm-based Linear Regression (WOA-LR), outperforms Genetic Algorithm-Linear Regression (GA-LR), Particle Swarm Optimization-Linear Regression (PSO-LR), JAYA-LR, and traditional Linear Regression (LR) in achieving desired objective functions and significantly reducing Mean Squared Error (MSE). This approach holds promise for more accurate resource utilization prediction and optimization in dynamic cloud environments.

*Keywords—Cloud computing; resource utilization; prediction; linear regression; metaheuristics*

## I. INTRODUCTION

With recent advances in artificial intelligence, the Internet of Things (IoT) [1, 2], Wireless Sensor Networks (WSNs) [3], and cloud computing, research efforts are shifting towards simplifying communication across various devices. Cloud computing represents a novel computational paradigm for provisioning computing resources, catering to a wide spectrum of users, ranging from individuals to large-scale enterprises [4, 5]. Cloud computing ecosystems are dominated by the intricate and cost-intensive data centers (DCs) that significantly impact service providers' financial viability [6]. Cloud providers offer three primary service categories, namely Software as a Service (SaaS), Platform as a Service (PaaS), and Infrastructure as a Service (IaaS), leveraging web service technology [7, 8]. Notable examples include Amazon for IaaS [9], Google for PaaS [10], and Salesforce for SaaS [11], all renowned as leading cloud providers worldwide. On one front, cloud providers offer their computational resources to fulfill users' Quality of Service (QoS) requirements, and on the other, they must effectively manage their Total Cost of Ownership (TCO) to thrive in the increasingly competitive cloud market [12]. Virtualization technology is widely employed within DCs to optimize resource allocation and reduce overall power consumption, a pivotal component of TCO. Furthermore, power management aligns with sustainability objectives. In a virtualized environment, a hypervisor intervenes to multiplex the resources of Physical Machines (PMs) among Virtual Machines (VMs) [13]. Inefficient resource allocation has repercussions on both resource utilization and the overall power consumption of DCs [14, 15].

Given the dynamic and ever-changing nature of cloud infrastructures and platforms, live VM migration emerges as a practical strategy that aligns with the current state of DCs [16]. Two common occurrences in DCs are under-utilization and over-utilization events [17]. The former entails a high-power consumption rate, while the latter is characterized by a high SLA violation rate. To address the scenario where businesses deploy VMs with varying usage patterns and resource requirements over time, machine learning approaches prove invaluable in discerning near-precise usage patterns in the short-term future [18]. Consequently, live virtual machine migration is employed to meet requirements before the aforementioned unfavorable events materialize. To maintain the desired QoS for users, Service Level Agreements (SLAs) are established between users and providers [19]. For instance, if a user submits an application comprising 150,000 million Instructions (MIs) and requests a VM with processing power equivalent to 250 million Instructions Per Second (MIPS), the provider must ensure the VM operates continuously at 100% utilization to deliver the results within 10 minutes. Failure to do so results in an SLA violation and a penalty for the service provider. The risk of SLA violation escalates when PMs become overloaded in DCs. Hence, preemptive offloading of some VMs prior to this event can mitigate SLA violations. Conversely, the phenomenon of server sprawl significantly increases total power consumption, whereby numerous under-loaded active PMs run concurrently. Server consolidation, which consolidates VMs into the fewest active PMs, is an effective strategy for reducing overall power consumption [20].

The integration of machine learning, deep learning, and meta-heuristic algorithms significantly enhances resource usage prediction in dynamic cloud environments. These methodologies provide a key insight into complex resource utilization patterns, contributing to the efficient management of cloud infrastructures [21]. Machine learning models, particularly linear regression and its variants, empower predictions based on historical resource usage data, enabling proactive resource allocation and load balancing [22]. Deep learning techniques, with their ability to analyze vast amounts of unstructured data, facilitate the identification of intricate patterns within cloud workloads, leading to more accurate predictions [23, 24]. Moreover, the inclusion of meta-heuristic

algorithms, such as genetic algorithms, particle swarm optimization, and whale optimization, augments predictive accuracy by refining the traditional models, enabling them to adapt and evolve according to dynamic resource demands [25]. The synergy among these methodologies results in robust and adaptable prediction models vital for optimizing resource usage in dynamic cloud environments, ultimately leading to improved service quality, cost efficiency, and better user experiences within cloud services [26, 27].

In this paper, machine learning techniques are extensively leveraged to derive resource usage patterns from historical data. This empowers the hypervisor to make swift decisions regarding under-utilization and over-utilization events before they occur. To address this challenge, machine learning techniques are applied to historical data to predict near-future resource requests. Processing and memory resources are pivotal for each requested VM, with processing capacity holding greater significance among power consumers, which informs the focus of this study on CPU capacity requirements [28]. To forecast near-future resource requests, the linear regression algorithm, a branch of machine learning, is employed. This involves recording the average resource utilization at five-minute intervals, utilizing the data history from the previous hour to inform short-term predictions. The time interval for data collection can be tailored to specific requirements. To enhance the performance of traditional linear regression, several effective meta-heuristic approaches are incorporated, yielding promising results. The innovation in this paper centers on the following key aspects:

- Introduction of a resource usage prediction model based on historical data.

- Presenting a migration trigger model for timely decision-making to avert unforeseen events.

- Exploration and evaluation of customized meta-heuristic algorithms to determine the most efficient approach.

The structure of this paper is organized as follows: Section II provides an overview of related work in the field. Section III introduces the proposed system model. In Section IV, the problem is formally defined and elaborated upon. Section V presents the suggested algorithm to address the problem. The performance assessment of the proposed approach is outlined in Section VI. Section VII concludes the paper and outlines potential future directions for research in this domain.

## II. BACKGROUNDS

Live VM migration is the intricate process of seamlessly transferring a VM's loads from one PM to another, ensuring uninterrupted service for the end user. The initiation of live VM migration and server consolidation is motivated by various factors, with power management being of paramount importance [29]. Other driving factors include mobile computing [30], reduction of communication costs [31, 32], system maintenance [33], and enhancing system failure reliability [34]. This leads to fundamental questions regarding when and where VM migration should be triggered, a topic that has been extensively explored in existing literature. One

notable contribution in this domain was made by Martinovic, et al. [35], who introduced a server consolidation model aimed at minimizing power consumption. They approached the problem by transforming the VM placement challenge into a bi-packing problem with conflicts and modeled by an integer linear program to solve it while utilizing the minimum number of PMs required. Zhou, et al. [36] proposed a linear regression model for predicting the CPU demands of VMs and subsequently triggering live VM migration in anticipation of near-future overload. Although this approach holds promise, it suffers from relatively high prediction errors.

Zhao, et al. [37] introduced a communication-aware live VM migration algorithm based on the Ant Colony Optimization (ACO) algorithm to minimize overall costs. This algorithm consists of two phases: first, it identifies VMs with high affinity for migration, and second, it selects the destination PMs for relocating the VMs. An energy-aware VM migration model was presented by Patel, et al. [38] to achieve both load balancing and power conservation, involving a three-way decision-making process for heavy, medium, and light workloads. A combined forecasting and load-aware migration model, along with an automated algorithm, was proposed by Forsman, et al. [39] to address live VM migration. A similar approach was put forth by Paulraj, et al. [40], focusing on saving energy, enhancing system reliability against failures, and maximizing service availability. The VM migration process, being resource-intensive and potentially degrading performance, employs forecasting models to estimate the resource requirements of each VM. Optimal online deterministic VM placement and adaptive heuristic algorithms are employed to address server consolidation, contributing to efficient DC power management and performance maintenance [41]. In the proposed approach, a novel system framework is introduced, featuring both local and global managers. The global manager resides at the master node, while each PM is equipped with a local manager responsible for collecting information on resource utilization and transmitting it to the global manager. Subsequently, the global manager issues directives for optimal VM placement, considering user SLAs. To further enhance efficiency, a cost function is defined, linked to the time associated with the migration process. The migration process is carefully orchestrated to avoid violating predetermined SLAs.

The review of existing literature reveals promising contributions from the research community. However, there remains a challenge in achieving near-precise prediction models. This motivates the current article, which leverages machine learning methods to predict short-term resource requirements for each VM, thereby triggering relevant VM migration processes to reduce both power consumption and SLA violation rates significantly.

## III. SYSTEM MODEL

This section introduces the proposed system model and its components, along with a schematic example to demonstrate how the model functions. The proposed system model, depicted in Fig. 1, consists of two main parts: the front end and the back end. In the front end, users request specific types and specifications of VMs encapsulated in SLA format. These

requests are forwarded to a broker module that possesses knowledge of the underlying infrastructure's capabilities. The requested resources for each VM are logged in a repository, accumulating as historical data. Subsequently, the forecaster module, a part of the live VM migration scheme, is activated based on the historical data in the repository. Its primary objective is to prevent both SLA violations resulting from overloaded conditions and high-power consumption due to the server sprawl phenomenon.

Fig. 2 provides a schematic example of a scenario involving the execution of different VMs within a data center with three PMs. Initially, three different users submit their requests denoted as $R_1$, $R_2$, and $R_3$, where $R_1$ entails a request for 2 VMs, $R_2$ for 1 VM, and $R_3$ for 2 VMs. The broker promptly dispatches these requests to the available PMs, as depicted in Fig. 2(a). At five-minute intervals, resource requests are recorded in a repository. The forecaster module

extracts insights from this data history and anticipates that $PM_3$ will become overloaded in the near short-term future due to the surging resource request of $VM_5$. To prevent an overload event, the live VM migration module is activated, and $PM_2$ is chosen for offloading due to its surplus resources. After the live migration, the scenario is represented in Fig. 2(b). In this situation, when the five-minute time interval is reached, the forecaster module predicts that $VM_2$ and $VM_4$ will reduce their resource requirements. This prediction aims to decrease resource requests.

Consequently, both $PM_1$ and $PM_3$ are in an under-loaded condition in the near future. Therefore, the live VM migration scheme is initiated for both $VM_2$ and $VM_4$ to consolidate servers. Subsequently, the unused $PM_3$ is transitioned into hibernation mode to conserve energy that would otherwise be dissipated.



Fig. 1.   System model.



Fig. 2.   A schematic example of a data center with deployed VMs.

## IV. PROBLEM STATEMENT

The core problem at hand is the necessity to make prompt decisions to preempt unfavorable events before they materialize. Periodically, the forecaster module retrieves data from the data repository within the data center to predict the impending over-loaded and under-loaded states of each PM. Subsequently, the most appropriate decision is made based on these predictions. To accomplish this, an advanced linear regression model is employed with the objective of minimizing the model's Mean Squared Error (MSE). In this pursuit, the conventional linear regression method is fused with a meta-heuristic algorithm, resulting in a novel and advanced hybrid forecaster algorithm. The key aim is to reduce the MSE to the lowest extent, thereby enhancing the accuracy of the decision-making process. In essence, the decision becomes increasingly accurate as the error approaches its minimum value. To facilitate this, the data history stored in the repository pertaining to resource requests within the most recent hour is divided into 12 records, each representing a five-minute interval. Utilizing the information within these recorded data, a linear function is established. Eq. (1) represents this linear function, where $x$ and $y$ denote the input and forecast functions, respectively. The terms $C_0$ and $C_1$ signify the two constants that serve as the coefficients of the linear function and must be determined by the proposed model.

$$y = C_0 + C_1 x \qquad (1)$$

## V. PROPOSED ALGORITHM

To address the live VM migration, which is inherently an optimization challenge, various competitive algorithms have been proposed. These algorithms have been selected based on their demonstrated success in the existing literature for solving continuous optimization problems and their adaptability to the specific problem under consideration. In this section, we introduce the following algorithms for calculating both LR's parameters $C_0$ and $C_1$: Canonical Linear Regression (LR), Genetic Algorithm-based Linear Regression (GA-LR), Particle Swarm Optimization-based Linear Regression (PSO-LR), Whale Optimization Algorithm-based Linear Regression (WOA-LR), and JAYA Linear Regression (JAYA-LR). GA-LR predicts the utilization of each server for the short-term future based on previously recorded mean server CPU utilization. In this context, single-point crossover and random gene mutation operators are employed. Additionally, a tournament algorithm is formulated. These algorithms are designed to optimize the LR's parameters $C_0$ and $C_1$, thereby enabling more accurate predictions and decision-making in the context of the problem at hand.

In the proposed GA-LR, random individuals are generated to represent both populations associated with the two constants, $C_0$ and $C_1$, which serve as the coefficients in the linear regression function. Each record for each coefficient is incorporated into the linear regression function, and the difference between this value and the actual value in the dataset is regarded as the fitness value. This optimization problem is a straightforward minimization problem, where the goal is to refine the coefficients over successive rounds. For encoding, each chromosome consists of two segments: the first part encodes an integer value, and the second part encodes a real

number. It's worth noting that binary genes are used in the encoding process. The Tournament selection procedure aims to increase the likelihood of selecting promising chromosomes. In this procedure, $K$ chromosomes are randomly chosen from the populations, and the best-performing ones are returned. It is important to highlight that the proposed tournament selection approach does not directly select the best individuals from the entire population, as this would risk early convergence, potentially resulting in suboptimal performance. The algorithm iterates until a specified termination condition is met, ultimately returning the best-performing chromosomes that yield the minimum MSE value.

In the PSO-LR algorithm, two distinct swarms of particles are randomly generated, akin to the populations of individuals in genetic Algorithms. Each particle's future trajectory is determined by three key parameters: inertia, local best, and global best values. The first parameter, inertia, is responsible for the particle continuing in its previous direction. The second parameter directs the particle to adjust its direction based on its local best, which is recorded in its memory. The third parameter steers the particle towards changing its direction to align with the global best of the entire swarm. To weigh the effectiveness of each parameter, specific weights are assigned to them. The algorithm is executed over multiple iterations, and subsequently, the best-performing particle thus far is identified and returned as the optimal solution. This iterative process helps refine the solution and converge towards the most accurate values for $C_0$ and $C_1$ in the linear regression model. For another comparative approach, the WOA-LR algorithm is presented in Fig. 3.

Similar to other swarm-based meta-heuristic algorithms, the WOA-LR begins with the generation of random swarms of whales. In this algorithm, each pair of whales represents a solution, as each solution requires two coefficients. To this end, variables $PopC_0(i)$ and $PopC_1(i)$ pertain to the $i^{th}$ whale. The loop encompassing lines 4 to 34 is executed for each whale, which is why there are two inner for-loops. In lines seven and eight, two sets of vectors, $a$, $A$, and $C$, are updated. It is important to note that $a$ is a vector that gradually decreases from 2 to 0. Additionally, the vectors $A$ and $C$ consist of random real values within the [0..1] interval. The changes in $r$ are determined by Eq. (2), and Eq. (3) defines the alterations in these vectors.

$$\vec{A} = 2\vec{a}.\vec{r} - \vec{a} \qquad (2)$$

$$\vec{C} = 2\vec{r} \qquad (3)$$

Furthermore, random real values $l$ are stochastically selected from the interval [-1..1]. The essence of the WOA lies in the oscillation between exploration and exploitation phases at intervals during the algorithm's lifecycle. To this end, a random variable $P$ is drawn to determine whether to explore or exploit the search space. The update process in line 15 is specifically designed to reflect the inclination towards exploitation or local search. In the context of exploration, the update is carried out using Eq. (4), and this operation is executed in line 18. In this scenario, the update is impartially performed based on the random position of the whale without any bias.

**Algorithm 4**. WOA-Regression

**Input:**

    Swarm1 , Swarm2 : Swarm;

    SwarmSize : integer;

    MAXITER : integer;

**Output:**

    Optimal C0, C1, and future server utilization

1. Initialize the two whale populations $PopC_0$ and $PopC_1$ including $PopC_0(i)$ and $PopC_1(i)$ (i=1,2,…,PopSize)
2. Calculate the fitness fucntion for all y(i)= PopC0[i]+ PopC1[i]*x(i)
3. let PopC0* and PopC1* to be the coefficient of the best whale recognized so far.
4. While iteration is not reach to *MaxIteration* Do
5.     for i=1 To *PopSize* Do
6.       for j=1 To *PopSize* Do
7.         Update vectors a0, A0, and C0 based on Eqs.(6-7) ; and random variable *l*
8.         Update vectors a1, A1, and C1 based on Eqs.(6-7) ; and random variable *l*
9.         Draw new real random number 0≤P≤1;
10.         if (P<0.5) then
11.           if ( |A| < 1) then
12.             Update the whlale PopC0[i] position based on Eq. (8)
13.           else if (|A| ≥ 1) then
14.             Select a random whale PopC0[j] from population
15.             Update the whale PopC0[i] position based on Eq.(9) incorporating PopC0[j]
16.           end-if
17.         else (* P ≥0.5 *)
18.           Update the whale PopC1[i] position based on Eq.(10)
19.         end-if
20.         Draw new real random number 0≤P≤1;
21.         if (P<0.5) then
22.           if ( |A| < 1) then
23.             Update the whlale PopC1[i] position based on Eq. (9)
24.           else if (|A| ≥ 1) then
25.             Select a random whale PopC1[j] from population
26.             Update the whale PopC1[i] position based on Eq.(11) incorporating PopC1[j]
27.           end-if
28.         else (* P ≥0.5 *)
29.           Update the whale PopC1[i] position based on Eq.(13)
30.         end-if
31.         Call Clamping PopC0[i] and PopC1[i] to correct the whale position if it returns infeasible solution
32.       end-for
33.     find current best so far from population and put it to PopC0* and PopC1*
34. end-while
35. return PopC0*+PopC1*.CurrentUtilization

Fig. 3. Algorithm WOA-LR.

$$\vec{W_i}(t+1) = \vec{W_j}(t) - \vec{A}.\vec{D} \qquad (4)$$

The algorithm uses Eq. (5) to update the chosen whale's location in order to imitate the distinctive circular movement of whales, which is sometimes referred to as a "spiral update position." This mechanism is designed to mimic the distinctive movement pattern of whales during the optimization process.

$$\vec{W_i}(t+1) = \vec{D'}.e^{bl}.Cos(2\pi) + \vec{W*}(t) \qquad (5)$$

After the update is executed for all whales, in case any whale's solution becomes infeasible, the Clamping function is invoked to adjust the encoded solution within the appropriate and feasible space. Line 31 represents the application of the

Clamping (.) function. The specific implementation of the Clamping function can vary based on the context. Ultimately, the best whale found so far, which represents an efficient solution, is returned as the final solution. Another successful optimization algorithm, which has been introduced recently, is the JAYA algorithm. The JAYA algorithm is specifically designed for addressing continuous optimization problems. Similar to other meta-heuristic algorithms, it commences with a limited number of randomly generated solutions. In each iteration, the best and worst solutions found so far are identified. The primary objective is to approach the best solution while distancing from the worst one. Throughout the evolution of each solution, if a new solution improves in terms

of fitness value, it is accepted; otherwise, the previous version is retained. Each solution is iteratively adjusted to converge toward the best solution gradually found thus far.

## VI. RESULTS

This section is devoted to the performance evaluation of the proposed prediction model. To assess its effectiveness, various meta-heuristic-based algorithms have been put forward, and a comparative study has been conducted [42-44]. In this context, three distinct sets of datasets representing data from the last hour have been randomly generated. Table I provides an overview of these datasets. All of the selected algorithms operate on the same datasets, ensuring a level playing field for fair competition in the evaluation process. The simulation results are presented in Table II, with the reported values for the successful WOA-Regression algorithm highlighted. The key to the success of WOA-Regression lies in its ability to strike a balance between the exploration and exploitation phases during the optimization process, effectively optimizing the search process.

One of the most critical concerns in energy-intensive data centers is power consumption. Furthermore, the rate of SLA violations significantly affects users' decisions when it comes to adopting cloud services. In practice, users tend to abandon unreliable cloud providers that cannot meet the agreed SLA terms. Thus, ensuring a high-quality experience for users is of paramount importance. To address this issue, a power consumption model with a linear relationship to CPU utilization is introduced in Eq. (6).

$$P_{Current}^{PM_j} = \lambda_j \times P_{Full}^{PM_j} + (1 - \lambda_j) \times P_{Full}^{PM_j} \times U_{CPU}^{PM_j} \qquad (6)$$

The term $\lambda_j$ is employed to denote that an idle machine consumes a certain percentage of power compared to a fully loaded machine. Various research studies, including the current paper, commonly use $\lambda_j$ as 70% of the power consumed by a fully utilized machine. To calculate the CPU utilization of a PM, the summation of the utilization of all co-hosted VMs processing requests is considered, which can be obtained using Eq. (7). Additionally, the memory utilization of each PM is determined by summing the requested memory of all co-hosted VMs, a calculation that can be performed using Eq. (8). These measurements are essential for assessing and managing the resource utilization of each PM in the data center.

$$U_{CPU}^{PM_j} = \sum_{i=1}^{n} U_{CPU}^{VM_i} . x_{ij} \qquad (7)$$

$$U_{CPU}^{PM_j} = \sum_{i=1}^{n} U_{Mem}^{VM_i} . x_{ij} \qquad (8)$$

The binary decision variable, denoted as $x_{ij}$, indicates whether a VM is placed on a PM or not. If a VM is placed on a PM, this variable is set to 1; otherwise, it is set to 0. Additionally, the terms $U_{CPU}^{VM_i}$ and $U_{Mem}^{VM_i}$ are used to represent the CPU and memory bandwidth requirements of a VM $i$. Similarly, the terms $U_{CPU}^{PM_j}$ and $U_{Mem}^{PM_j}$ are employed to denote the CPU and memory utilization of a PM $j$, respectively. These variables and terms play a vital role in optimizing the allocation and utilization of resources within the data center.

TABLE I. AN OVERVIEW OF THE GENERATED DATASETS

| Server | CPU utilization recorded for different rounds | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | *Round 1* | *Round 2* | *Round 3* | *Round 4* | *Round 5* | *Round 6* | *Round 7* | *Round 8* | *Round 9* | *Round 10* | *Round 11* | *Round 12* |
| First server | 82% | 93% | 73% | 69% | 33% | 81% | 94% | 83% | 39% | 64% | 27% | 82% |
| Second server | 69% | 72% | 37% | 58% | 17% | 6% | 62% | 41% | 40% | 69% | 53% | 55% |
| Third server | 66% | 96% | 49% | 43% | 81% | 71% | 6% | 42% | 17% | 94% | 41% | 76% |

TABLE II. SIMULATION RESULTS

| Error calculation model | First PM | Second PM | Third PM |
|---|---|---|---|
| WOA-LR | MSE= 0.514 | MSE= 0.352 | MSE= 1.082 |
| Regression coefficient | C0= 0.594, C1= -0.093 | C0= 0.664, C1= -0.417 | C0= -7.601, C1= 1.271 |
| JAYA-LR | MSE= 1.128 | MSE= 0.505 | MSE= 1.282 |
| Regression coefficient | C0= 0.628, C1= -0.16 | C0= 0.784, C1= -0.511 | C0= 0.654, C1= -0.191 |
| PSO-LR | MSE= 1.0073 | MSE= 0.441 | MSE= 1.123 |
| Regression coefficient | C0= 0.624, C1= -0.752 | C0= 0.355, C1= 0.259 | C0= 0.751, C1= -0.341 |
| GA-LR | MSE= 0.684 | MSE= 0.432 | MSE= 1.108 |
| Regression coefficient | C0= 0.722, C1= -0.0481 | C0= 0.407, C1= 0.126 | C0= 0.563, C1= -0.116 |
| Conventional LR | MSE= 5.762 | MSE= 4.128 | MSE=7.361 |
| Regression coefficient | C0= 0.653, C1= 0.024 | C0= 0.397, C1= 0.154 | C0= 0.622, C1= -0.164 |

The live VM migration technique allows the transfer of a VM's pages from the source PM to the destination PM with minimal interruption, resulting in a small downtime. However, VM live migration can have adverse effects on the overall system performance and potentially jeopardize SLA between users and service providers. This is due to the time required for page transfers during migration. For cloud-based web applications, on average, downtime may decrease CPU utilization by approximately 10%. In other words, this phenomenon can lead to SLA violations to some extent. On the flip side, it is desirable to reduce the number of aggressive live migrations. The Live Migration Time (LMT) is highly dependent on the size of the VM's pages being transferred and the underlying bandwidth capacity. Since data centers often use Storage Area Networks (SAN), there is no need to transfer VM storage data, as all PMs have uniform access to SAN. The LMT for a VM is calculated using Eq. (9). This information is crucial for managing VM migrations efficiently and minimizing potential SLA violations.

$$LMT(VM_i) = \frac{Msize(M_i)}{BW_i} \qquad (9)$$

The term $Msize(M_i)$ represents the memory size of data being transferred via a shared link with a bandwidth capacity of $BW_i$. Furthermore, experimental results indicate that there is a 10% degradation in CPU utilization during the process of live migration. This performance degradation is quantified using Eq. (10), where the term $u_i(t)$ represents CPU utilization associated with VM $i$ during the migration process. This equation provides a measure of the performance impact of live migrations, which is crucial for optimizing resource allocation and minimizing SLA violations.

$$PD(VM_i = 0.1. \int_{t_0}^{t_0+LMT(VM_i)} u_i(t)dt \qquad (10)$$

The paramount issue that encourages users to remain loyal to specific cloud providers is the delivery of a high-quality user experience from the services provided. In this context, key points such as the minimum throughput and the maximum response time must be determined to meet the required QoS, which is outlined in the SLA. Web applications that leverage cloud infrastructure often exhibit fluctuations in resource utilization. As a result, an independent parameter reflecting the system's SLA violation rate is needed. To address this, two new parameters have been introduced: SLA violation length per active PM, denoted as $a$, and the total performance degradation due to VM migration, denoted as $\beta$. The parameter $\alpha$ represents the duration during which active PMs experience 100% CPU utilization, indicating the time span during which PMs are overloaded. This parameter is measured using Eq. (11) and plays a crucial role in quantifying SLA violations within the system.

$$a = \frac{1}{n}\sum_{i=1}^{n} \frac{T_{PM_i}}{T_{active(PM_i)}} \qquad (11)$$

where, $n$ represents the number of PMs, $T_{PM_i}$ stands for the time during which $PM_i$ experiences 100% CPU utilization, indicating when $PM_i$ is overloaded. $T_{active(PM_i)}$ is the total time during which $PM_i$ remains active, serving various VMs. The parameter $\beta$ is calculated using Eq. (12). In this equation, $m$ is the number of VMs, $PD(VM_i)$ is the performance degradation caused by VM migration, and $C(VM_i)$ represents the total CPU resource capacity associated with $VM_i$ in terms of MIPS.

$$\beta = \frac{1}{m}\sum_{i=1}^{m} \frac{PD(VM_i)}{C(VM_i)} \qquad (12)$$

Both parameters $\alpha$ and $\beta$ independently influence SLA violations. To provide a comprehensive assessment of SLA violations, a new parameter called the SLA Violation Rate (SLAVR) is introduced in Eq. (13).

$$SLAVR = \alpha.\beta = \frac{1}{n}\sum_{i=1}^{n} \frac{T_{PM_i}}{T_{active(PM_i)}} . \frac{1}{m}\sum_{i=1}^{m} \frac{PD(VM_i)}{C(VM_i)} \qquad (13)$$

With the inclusion of live migration costs, the performance of comparative algorithms is assessed in terms of energy consumption attributable to SLA violations, SLAVR, $\alpha$ (SLA violation length per active PM), $\beta$ (total performance degradation due to VM migrations), and the number of VM migrations. Table III provides a comparison of the state-of-the-art algorithms based on these assessment metrics. All of the comparative algorithms were executed in the CloudSim environment, with 20 independent runs. The reported results represent the average outcomes obtained from these runs.

TABLE III. PERFORMANCE EVALUATION

| Algorithm | VM migrations | β (%) | α (%) | SLAVR (%) | Total power consumption (Watt) | Power for SLAVR (Watt) |
|---|---|---|---|---|---|---|
| WOA-LR | 189 | 0.11 | 1.31 | 14.76 | 1327 | 65.11 |
| JAYA-LR | 203 | 0.13 | 1.36 | 15.61 | 1672 | 77.05 |
| PSO-LR | 208 | 0.13 | 1.37 | 17.79 | 1463 | 79.19 |
| GA-LR | 211 | 0.14 | 1.48 | 19.62 | 1495 | 75.91 |
| Conventional LR | 306 | 0.16 | 1.42 | 21.47 | 1588 | 89.02 |

## VII. CONCLUSION

This paper introduced a system framework for cloud data centers, comprising multiple modules designed to enable timely live VM migration for preventing SLA violations through the integration of machine learning tools. In this framework, a repository module is deployed within the DC, functioning as a data history repository where each physical machine records its average CPU utilization over time. Subsequently, a machine learning tool, specifically a linear regression-based model, is employed at regular intervals to predict near-future resource requirements. Based on these predictions, decisions are made to optimize resource allocation and avoid SLA violations.

There are still unanswered issues and unresolved problems in this field. Expanding the predictive capabilities to include important resources such as memory, storage, and network bandwidth might improve the overall effectiveness of the system. Furthermore, conducting performance evaluations on a range of real-world datasets and in different workload conditions would provide a more thorough understanding of the system's capacity to adjust and withstand challenges. Furthermore, it is necessary to do further research to determine the scalability, flexibility, and real-time responsiveness of the system when implemented and deployed in live cloud settings. The limitations of this study are its primary emphasis on CPU use, which may result in neglecting the complex interaction between various resources and their effect on adhering to SLA requirements. Moreover, the system's reliance on past data may provide difficulties in dynamic settings, requiring ongoing model training and adjustment methods. Future improvements may include integrating sophisticated machine learning methods, such as deep learning algorithms, to boost the accuracy of predictions and consider intricate resource linkages. Furthermore, investigating decentralized or distributed decision-making models for VM migrations, while also taking into account security and privacy concerns in a multi-tenant cloud environment, presents a promising direction for future study and improvement of the suggested framework.

## REFERENCES

[1] B. Pourghebleh, V. Hayyolalam, and A. A. Anvigh, "Service discovery in the Internet of Things: review of current trends and research challenges," Wireless Networks, vol. 26, no. 7, pp. 5371-5391, 2020.

[2] B. Pourghebleh and N. J. Navimipour, "Data aggregation mechanisms in the Internet of things: A systematic review of the literature and recommendations for future research," Journal of Network and Computer Applications, vol. 97, pp. 23-34, 2017.

[3] J. Zandi, A. N. Afooshteh, and M. Ghassemian, "Implementation and analysis of a novel low power and portable energy measurement tool for wireless sensor nodes," in Electrical Engineering (ICEE), Iranian Conference on, 2018: IEEE, pp. 1517-1522, doi: 10.1109/ICEE.2018.8472439.

[4] S. Vinoth, H. L. Vemula, B. Haralayya, P. Mamgain, M. F. Hasan, and M. Naved, "Application of cloud computing in banking and e-commerce and related security threats," Materials Today: Proceedings, vol. 51, pp. 2172-2175, 2022.

[5] V. Hayyolalam, B. Pourghebleh, M. R. Chehrehzad, and A. A. Pourhaji Kazem, "Single - objective service composition methods in cloud manufacturing systems: Recent techniques, classification, and future trends," Concurrency and Computation: Practice and Experience, vol. 34, no. 5, p. e6698, 2022.

[6] A. H. A. Al-Jumaili, R. C. Muniyandi, M. K. Hasan, J. K. S. Paw, and M. J. Singh, "Big Data Analytics Using Cloud Computing Based Frameworks for Power Management Systems: Status, Constraints, and Future Recommendations," Sensors, vol. 23, no. 6, p. 2952, 2023.

[7] S. AlMuraytib, L. Alqurashi, and S. Snoussi, "Blockchain-based solutions for Cloud Computing Security: A Survey," in Proceedings of the 6th International Conference on Future Networks & Distributed Systems, 2022, pp. 338-342.

[8] B. Kruekaew and W. Kimpan, "Multi-objective task scheduling optimization for load balancing in cloud computing environment using hybrid artificial bee colony algorithm with reinforcement learning," IEEE Access, vol. 10, pp. 17803-17818, 2022.

[9] X. Song, L. Pan, and S. Liu, "An online algorithm for optimally releasing multiple on-demand instances in IaaS clouds," Future Generation Computer Systems, vol. 136, pp. 311-321, 2022.

[10] R. Kaviarasan, P. Harikrishna, and A. Arulmurugan, "Load balancing in cloud environment using enhanced migration and adjustment operator based monarch butterfly optimization," Advances in Engineering Software, vol. 169, p. 103128, 2022.

[11] E. Ahumada-Tello and R. Evans, "A Complexity-based Framework for Social Product Development," Procedia CIRP, vol. 119, pp. 1204-1209, 2023.

[12] S. S. Gill et al., "Transformative effects of IoT, Blockchain and Artificial Intelligence on cloud computing: Evolution, vision, trends and open challenges," Internet of Things, vol. 8, p. 100118, 2019.

[13] H. Vahideh, P. Behrouz, P. K. A. Asghar, and A. Ghaffari, "Exploring the state-of-the-art service composition approaches in cloud manufacturing systems to enhance upcoming techniques," The International Journal of Advanced Manufacturing Technology, vol. 105, no. 1-4, pp. 471-498, 2019.

[14] K. Saidi and D. Bardou, "Task scheduling and VM placement to resource allocation in Cloud computing: challenges and opportunities," Cluster Computing, vol. 26, no. 5, pp. 3069-3087, 2023.

[15] S. Pazouki and M. R. Haghifam, "Optimal planning and scheduling of smart homes' energy hubs," International Transactions on Electrical Energy Systems, vol. 31, no. 9, p. e12986, 2021.

[16] P. Verma et al., "Voltage Rise Mitigation in PV Rich LV Distribution Networks Using DC/DC Converter Level Active Power Curtailment Method," Energies, vol. 15, no. 16, p. 5901, 2022.

[17] S. Durairaj and R. Sridhar, "MOM-VMP: multi-objective mayfly optimization algorithm for VM placement supported by principal component analysis (PCA) in cloud data center," Cluster Computing, pp. 1-19, 2023.

[18] A. Belgacem, S. Mahmoudi, and M. A. Ferrag, "A machine learning model for improving virtual machine migration in cloud computing," The Journal of Supercomputing, pp. 1-23, 2023.

[19] J. Singh and M. S. Goraya, "An Autonomous Multi-Agent Framework using Quality of Service to prevent Service Level Agreement Violations in Cloud Environment," International Journal of Advanced Computer Science and Applications, vol. 14, no. 3, 2023.

[20] B. Pourghebleh, A. A. Anvigh, A. R. Ramtin, and B. Mohammadi, "The importance of nature-inspired meta-heuristic algorithms for solving virtual machine consolidation problem in cloud environments," Cluster Computing, pp. 1-24, 2021.

[21] S. Vairachilai, A. Bostani, A. Mehbodniya, J. L. Webber, O. Hemakesavulu, and P. Vijayakumar, "Body Sensor 5 G Networks Utilising Deep Learning Architectures for Emotion Detection Based On EEG Signal Processing," Optik, p. 170469, 2022.

[22] S. P. Rajput et al., "Using machine learning architecture to optimize and model the treatment process for saline water level analysis," Journal of Water Reuse and Desalination, 2022.

[23] V. Monjezi, A. Trivedi, G. Tan, and S. Tizpaz-Niari, "Information-Theoretic Testing and Debugging of Fairness Defects in Deep Neural Networks," presented at the 2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE), 2023. [Online]. Available: https://doi.ieeecomputersociety.org/10.1109/ICSE48619.2023.00136.

[24] W. Anupong et al., "Deep learning algorithms were used to generate photovoltaic renewable energy in saline water analysis via an oxidation process," Water Reuse, vol. 13, no. 1, pp. 68-81, 2023.

[25] S. R. Abdul Samad et al., "Analysis of the Performance Impact of Fine-Tuned Machine Learning Model for Phishing URL Detection," Electronics, vol. 12, no. 7, p. 1642, 2023.

[26] M. Hajihosseinlou, A. Maghsoudi, and R. Ghezelbash, "Stacking: A novel data-driven ensemble machine learning strategy for prediction and mapping of Pb-Zn prospectivity in Varcheh district, west Iran," Expert Systems with Applications, vol. 237, p. 121668, 2024.

[27] M. Hajihosseinlou, A. Maghsoudi, and R. Ghezelbash, "A Novel Scheme for Mapping of MVT-Type Pb–Zn Prospectivity: LightGBM, a Highly Efficient Gradient Boosting Decision Tree Machine Learning Algorithm," Natural Resources Research, pp. 1-22, 2023.

[28] M. H. Shirvani, "An energy-efficient topology-aware virtual machine placement in Cloud Datacenters: A multi-objective discrete JAYA optimization," Sustainable Computing: Informatics and Systems, vol. 38, p. 100856, 2023.

[29] K. Kumar, K. Patange, P. Pete, M. Wankhade, A. Chatterjee, and M. Kurhekar, "Power and Energy-efficient VM scheduling in OpenStack Cloud Through Migration and Consolidation using Wake-on-LAN," IETE Journal of Research, pp. 1-13, 2022.

[30] Y. Kumar, S. Kaul, and Y.-C. Hu, "Machine learning for energy-resource allocation, workflow scheduling and live migration in cloud computing: State-of-the-art survey," Sustainable Computing: Informatics and Systems, vol. 36, p. 100780, 2022.

[31] S. Manjunatha and L. Suresh, "Optimal Min-Communication and Migration Cost Algorithm based Approach for Efficient Task Migration in Cloud Computing," in 2021 International Conference on Circuits, Controls and Communications (CCUBE), 2021: IEEE, pp. 1-6.

[32] S. Shahryari, F. Tashtarian, and S.-A. Hosseini-Seno, "CoPaM: Cost-aware VM Placement and Migration for Mobile services in Multi-Cloudlet environment: An SDN-based approach," Computer Communications, vol. 191, pp. 257-273, 2022.

[33] R. M. Haris, K. M. Khan, and A. Nhlabatsi, "Live migration of virtual machine memory content in networked systems," Computer Networks, vol. 209, p. 108898, 2022.

[34] M. Torquato, P. Maciel, and M. Vieira, "Availability and reliability modeling of vm migration as rejuvenation on a system under varying workload," Software Quality Journal, vol. 28, pp. 59-83, 2020.

[35] J. Martinovic, M. Hähnel, G. Scheithauer, and W. Dargie, "An introduction to stochastic bin packing-based server consolidation with conflicts," Top, vol. 30, no. 2, pp. 296-331, 2022.

[36] H. Zhou, Q. Li, K.-K. R. Choo, and H. Zhu, "DADTA: A novel adaptive strategy for energy and performance efficient virtual machine consolidation," Journal of Parallel and Distributed Computing, vol. 121, pp. 15-26, 2018.

[37] H. Zhao et al., "VM performance-aware virtual machine migration method based on ant colony optimization in cloud environment," Journal of Parallel and Distributed Computing, vol. 176, pp. 17-27, 2023.

[38] D. Patel, R. K. Gupta, and R. Pateriya, "Energy-aware prediction-based load balancing approach with VM migration for the cloud environment," Data, Engineering and Applications: Volume 2, pp. 59-74, 2019.

[39] M. Forsman, A. Glad, L. Lundberg, and D. Ilie, "Algorithms for automated live migration of virtual machines," Journal of Systems and Software, vol. 101, pp. 110-126, 2015.

[40] G. J. L. Paulraj, S. A. J. Francis, J. D. Peter, and I. J. Jebadurai, "A combined forecast-based virtual machine migration in cloud data centers," Computers & Electrical Engineering, vol. 69, pp. 287-300, 2018.

[41] A. Beloglazov and R. Buyya, "Optimal online deterministic algorithms and adaptive heuristics for energy and performance efficient dynamic consolidation of virtual machines in cloud data centers," Concurrency and Computation: Practice and Experience, vol. 24, no. 13, pp. 1397-1420, 2012.

[42] K. Ramana, R. Aluvalu, V. K. Gunjan, N. Singh, and M. N. Prasadhu, "Multipath Transmission Control Protocol for Live Virtual Machine Migration in the Cloud Environment," Wireless Communications and Mobile Computing, vol. 2022, 2022.

[43] A. Gupta, P. Dimri, and R. Bhatt, "An Optimized Approach for Virtual Machine Live Migration in Cloud Computing Environment," in Evolutionary Computing and Mobile Sustainable Networks: Springer, 2021, pp. 559-568.

[44] K. J. Naik, "An Adaptive Push-Pull for Disseminating Dynamic Workload and Virtual Machine Live Migration in Cloud Computing," International Journal of Grid and High Performance Computing (IJGHPC), vol. 14, no. 1, pp. 1-25, 2022.

# Efficient Processing of Large-Scale Medical Data in IoT: A Hybrid Hadoop-Spark Approach for Health Status Prediction

Yu Lina[1]*, Su Wenlong[2]

Hebei College of Industry and Technology, Hebei Shijiazhuang, 050091, China[1]
Liaoning University, Liaoning Shenyang, 110036, China[2]

*Abstract*—In the realm of Internet of Things (IoT)-driven healthcare, diverse technologies, including wearable medical devices, mobile applications, and cloud-based health systems, generate substantial data streams, posing challenges in real-time operations, especially during emergencies. This study recommends a hybrid architecture utilizing Hadoop for real-time processing of extensive medical data within the IoT framework. By employing distributed machine learning models, the system analyzes health-related data streams ingested into Spark streams via Kafka threads, aiming to transform conventional machine learning methodologies within Spark's real-time processing, crafting scalable and efficient distributed approaches for predicting health statuses related to diabetes and heart disease while navigating the landscape of big data. Furthermore, the system provides real-time health status forecasts based on a multitude of input features, disseminates alert messages to caregivers, and stores this valuable information within a distributed database, which is instrumental in health data analysis and the production of flow reports. We compute a range of evaluation parameters to evaluate the proposed methods' efficacy. This assessment phase encompasses measuring the performance of the Spark-based machine learning algorithm in a distributed parallel computing environment.

*Keywords—Internet of Things; big data; hadoop; spark-based machine learning*

## I. INTRODUCTION

Over the past two decades, our epoch has come to be recognized as the era of big data, wherein digital data has assumed a pivotal role across various domains, encompassing society, research endeavors, and, particularly, the medical domain [1]. Big data denotes the characterization of copious data amassed from diverse sources, such as sensor networks, high-throughput apparatus, mobile applications, streaming devices, and data reservoirs spanning numerous industries, with a pronounced emphasis on the healthcare sector [2, 3]. Effectively managing, processing, presenting, and deriving insights from this diverse and voluminous data spectrum has posed substantial challenges using the extant technological toolset [4]. Efficiently deriving meaningful insights from this multitude of data, tailored to various user profiles, ranks among the paramount technological quandaries facing the domain of big data analytics [5, 6]. Presently, numerous data sources within healthcare, both clinical and non-clinical, are converging, with the digital medical history of patients being of paramount importance in healthcare analytics [7].

Consequently, three primary challenges surface in creating a distributed data system designed to handle extensive data volumes [8].

The initial challenge stems from the complexity of collecting data from disparate sources due to its heterogeneous and vast nature. Second, the fundamental predicament revolves around storage, as big data systems must effectively store data while maintaining optimal performance. The final challenge pertains to big data analytics, especially real-time or near-real-time analysis of vast datasets, incorporating forecasting, optimization, visualization, and modeling [9]. In light of the shortcomings of current data management systems in addressing real-time and heterogeneous data, a need emerges for a new processing paradigm [10]. Conventional relational database management systems, exemplified by MySQL, predominantly cater to structured data management, with limited support for unstructured or partially structured data. Furthermore, traditional RDBMS scaling strategies for parallel hardware management and fault tolerance often prove inadequate as data volumes expand [11].

To tackle these challenges, the research community has introduced a variety of projects to address large-scale and diverse data storage, including NoSQL database management systems suitable for scenarios where a relational model is not requisite. MapReduce, an amalgamation of Map and Reduce operations, serves as a parallel processing technique for handling vast distributed datasets in commodity clusters [12]. Yet, it is marred by its sluggishness when dealing with iterative algorithms. The Hadoop framework, a batch processing system, is employed for distributed data processing and storage, relying on the MapReduce model for programming [13]. The Hadoop Distributed File System (HDFS) offers a distributed storage solution that is highly resilient [14]. However, Hadoop is ill-suited for in-memory computing and real-time stream processing and does not uniformly apply the MapReduce paradigm to all challenges. The volume of processed data is a determinant of the speed of results. Conversely, stream computing prioritizes data velocity and involves continuous input and output. Big data streaming computing (BDSC) comprises real-time computing, distributed messaging, high throughput, and minimized processing latency. It is essential for extracting meaningful information from vast datasets, particularly in the healthcare realm.

The swift advancement of large data analytics holds significant implications for advancing medical practices and academic research. Data collection, management, analysis, and assimilation tools designed to handle heterogeneous, unstructured, and structured data within contemporary healthcare systems have become accessible. BDSC is now integral to the landscape of big data analytics, facilitating the rapid exploration of the latent value of extensive healthcare data. Nevertheless, challenges persist due to the diverse data sources within the healthcare sector, necessitating the integration of data originating from relational databases, Hadoop, search engines, and other analytical systems. The application of machine learning to such extensive and high-velocity data streams presents considerable challenges, as conventional machine learning algorithms are not well-suited for such massive data volumes and variable velocities. Furthermore, efficient analytical data processing is a pressing concern, necessitating effective data integration. While contemporary research predominantly relies on machine learning, real-time machine learning applications are absent for streaming big data. Moreover, most healthcare analytics solutions predominantly focus on Hadoop, a batch-oriented computational platform.

The growing elderly population and the rising prevalence of chronic illnesses have exacerbated the inadequacies of conventional healthcare practices. In tandem, medical IoT has increased, enabling continuous monitoring and real-time emergency interventions, especially in cardiac conditions. This proliferation has led to the generation of vast datasets by millions of sensors, challenging the capacity to process and respond to this data under critical conditions. To address these challenges, we have developed a healthcare framework exemplified by a real-time health status forecasting case study. NoSQL Cassandra, Spark streaming, Spark MLlib, Kafka data streaming, and Apache Zeppelin technologies underpin this system. Kafka's producers generate multiple message streams, which are filtered using Spark streaming, enriched through machine learning, and stored in NoSQL repositories, facilitating analytics and visualization. This endeavor has substantially improved the quality of patient monitoring within healthcare.

The remaining portion of the paper is organized in the following fashion. Section II provides a comprehensive analysis of previous research in the field, which serves as a foundation for our suggested approach. Section III provides a detailed explanation of the hybrid architecture, with a focus on the incorporation of Hadoop, Spark, and distributed machine learning models. Section IV provides detailed explanations of the specific scenarios or use cases relevant to our proposed architecture. Subsequently, Section V provides the discussion of comprehensive examination, evaluating the architecture's advantages, constraints, possible uses, and comparative observations. Section VI explores the collected data, demonstrating the effects of adopting our architecture for predicting health status. Section VII ultimately ends by providing a concise overview of significant discoveries and proposing potential avenues for further study.

## II. RELATED WORKS

### A. Medical Big Data Challenges

The concept of the 5Vs in big data, comprising Volume, Variety, Velocity, Veracity, and Value, aptly elucidates the sheer magnitude of data generated within the contemporary healthcare sector. The healthcare domain is burdened by a substantial and ever-expanding volume of data that necessitates comprehensive collection and analysis [15]. The notion of variety underscores the diverse range of data sources that must be tapped into within healthcare. Pertinently, healthcare data and the domain's knowledge demand real-time acquisition, encapsulated by the concept of velocity. The integrity and trustworthiness of healthcare data are encapsulated in the dimension of veracity. Ultimately, valuable insights can be gleaned through meticulous examination of the colossal healthcare dataset. Distributed sources of healthcare data encompass medical electronic records, health claims, diagnosis data, clinical imagery, streaming systems, and sensors affixed to patients' bedsides for continuous vital sign monitoring. These sources collectively generate vast amounts of data, surpassing the processing capabilities of conventional data handling systems. The myriad challenges associated with big data are illustrated in Fig. 1. In this research, our focus has been dedicated to the initial five pivotal challenges within big data, encompassing data integration, storage, analysis, and representation.



Fig. 1. Big data challenges.

### B. Literature Study

The amalgamation of Machine Learning (ML), Deep Learning (DL), Neural Networks (NN), Fuzzy Logic Systems (FLSs), Wireless Sensor Networks (WSNs), and Temporal Graphs (TGs) holds pivotal significance in the processing of large-scale medical data within the IoT landscape. ML and DL techniques empower healthcare systems to discern intricate patterns within vast datasets, enabling predictive analytics for disease diagnosis, treatment planning, and health status forecasting [16-18]. NN, a subset of ML, simulates the human brain's learning process, aiding in complex data analysis, especially in image recognition and signal processing tasks

within medical imaging and diagnostics [19, 20]. FLSs supplements decision-making processes by handling uncertain or imprecise data, crucial in medical scenarios where data might exhibit variability [21]. WSNs, integrated with IoT devices, facilitate real-time health monitoring, efficiently collecting and transmitting patient data for timely analysis [22]. Meanwhile, TGs provide an intricate understanding of dynamic patient interactions over time, aiding in disease progression modeling and personalized treatment plans [23]. The convergence of these technologies optimizes medical data processing, fostering precision medicine, remote patient monitoring, and efficient healthcare delivery, thereby revolutionizing patient care and augmenting medical research endeavors within the IoT-driven healthcare domain [24].

The exponential proliferation of healthcare data, coupled with its profound insights, has positioned big data analytics, particularly within the healthcare domain, as a formidable challenge spanning multiple academic disciplines, including data mining and machine learning. The advancement of data collection in healthcare can be primarily attributed to strides in scientific and technological innovation. The healthcare sector primarily leverages three fundamental categories of digital data for data collection: health research, operational processes within healthcare organizations, and clinical records. Traditional data mining techniques, which involve identifying valuable patterns within vast databases, struggle to unearth insights from the dispersed, extensive, and diverse datasets prevalent in healthcare. Data mining techniques are pivotal in transforming this data into actionable information. Numerous studies in the medical field focus on prediction and recommendation systems. These studies include experiments on heart attack prediction and a comparative assessment of different approaches. Breast carcinoma classification employs a genetically tuned neural network model. Other research endeavors encompass information retrieval and data mining methodologies.

Healthcare analytics encompass various applications, such as epidemic forecasting, health decision support, and recommendation systems, geared toward enhancing care quality, reducing costs, and augmenting productivity. One notable approach is utilizing a K-means clustering algorithm operating in the cloud as a MapReduce task, utilizing healthcare data for clustering. An alternative proposal suggests a decentralized platform for managing electronic health records personally, employing Hadoop and HBase. Predictive analysis in healthcare involves forecasting diabetes and determining the most suitable therapy using algorithms and the Hadoop MapReduce environment. Big data contextual exchange among healthcare systems through the Internet of Things (IoT) is demonstrated through an intelligent care system built on Hadoop. This system leverages an architecture with advanced data processing capabilities to collect data from diverse linked devices and transmit it to intelligent buildings. Real-time analysis of electronically generated medical records and data from medical equipment and mobile applications is described. This system, incorporating Hadoop, MongoDB, and an innovative treatment method, aims to enhance patient information processing outcomes. The predominant focus in most healthcare analytics solutions lies in Hadoop, which can

handle substantial data volumes from diverse sources in batch-oriented processing. However, Hadoop's real-time processing capabilities are limited, and Spark emerges as a swifter and more efficient alternative, particularly for iterative machine learning tasks. Both Hadoop and Spark, being Apache projects, are integral to the big data landscape, with Spark generating significant interest.

Several scalable machine learning algorithms aim to address the diverse challenges within big data analytics. These algorithms include a scalable Random Forest classification model for diabetes risk prediction, logistic regression for phishing URL detection, and a Markov chain-based system for identifying abnormal patterns in the behaviors of elderly individuals. Real-time management of medical emergencies using IoT-based medical sensors is presented, along with a paradigm for real-time analysis of extensive medical data using Spark Streaming and Apache Kafka. A real-time health forecasting system focusing on machine learning, particularly Decision Trees, is developed to process data streams obtained via socket streams. A novel strategy for cardiac disease monitoring, centered on real-time decentralized machine learning within the Spark environment, is proposed in one study. Most of these studies either center on specific healthcare data sources or predominantly deal with batch-oriented computation. Healthcare generates a myriad of rapidly accumulating data from diverse sources. Moreover, some studies prioritize data storage and visualization, while others emphasize powerful data analytics tools like data mining and machine learning. Thus, the creation of an effective system for managing remote health data streams necessitates real-time healthcare analysis, which encompasses data collection, real-time processing, and robust machine learning capabilities.

The two leading causes of global mortality in recent times have been heart disease and diabetes. Continuous monitoring and early detection of these ailments can significantly reduce mortality rates. The availability of wearable health monitors, the adoption of IoT medical technology within healthcare systems, and the surge in patient conditions further underscore the potential of big data technologies for real-time health condition prediction. Real-time prediction can streamline healthcare visits and empower patients and healthcare providers to anticipate potential illnesses. Furthermore, the proposed system includes an alert mechanism, ensuring that emergency services are promptly notified when a patient's condition deviates from the norm, facilitating rapid interventions during emergencies.

### III. PROPOSED ARCHITECTURE

Within the scope of this research, a system for data processing and monitoring is introduced, amalgamating Kafka and Spark streams. This system operates by first processing data received from connected devices and subsequently storing this data for real-time analysis. The architectural layout of this proposed system is elucidated in Fig. 2. The system commences with the continuous generation of data messages from Kafka generators. These data messages encompass diverse disease names and are subsequently conveyed to a Spark streaming application for immediate processing. Spark Stream harnesses machine learning models to analyze various

health attributes acquired from the Kafka Stream, thereby predicting health status. The results of this analysis are stored in a NoSQL Cassandra database. In the proposed architecture, Apache Zeppelin is instrumental in retrieving data from the database and presenting it in a real-time dashboard featuring data visualization in the form of graphs, charts, and data tables. By leveraging this real-time data in an Internet of Things (IoT)

context, it becomes feasible to promptly scrutinize it, enabling the timely dispatch of alert notifications to caregivers when significant changes in a patient's condition arise. This real-time monitoring capability facilitates immediate action and intervention when necessary, ensuring patients receive timely and responsive care.



Fig. 2. Proposed architecture.

The IoT encompasses a network of physical and virtual entities equipped with electronics, intelligent wearables, software, applications, sensors, and network connectivity, all designed to collect and exchange data among themselves and with data center systems. The data generated by ubiquitous wearable health monitors, commonly found in households, is characterized by its substantial volume and random nature. Stimulating user activity trends or gathering essential data necessitates analysis through a robust big data analytics system. Forecasts indicate that, by 2020, IoT-related technologies within the healthcare sector will constitute a significant portion, accounting for forty percent of all IoT-related technologies. The integration of information technology in healthcare, particularly health informatics, is poised to bring about a paradigm shift, significantly reducing inefficiencies, containing costs, and, ultimately, saving lives. Real-time monitoring facilitated by the IoT can be a lifesaver in medical emergencies, encompassing conditions such as diabetes, heart disease, and various chronic disorders. Numerous sources are presently accessible for the continuous monitoring of health indicators. The workflow of the proposed system, involving multiple data sources, is outlined in Fig. 3. This system aims to harness the power of IoT to provide real-time monitoring and timely interventions in healthcare, thereby enhancing the quality of care and potentially saving lives in critical situations.

The escalating volume of data generated within healthcare systems has surpassed the capabilities of Spark alone for data management. In response to this challenge, Kafka, designed explicitly for managing streaming data, has been seamlessly integrated into our system. The data collection component in the proposed system architecture plays a pivotal role in gathering health-related data from various sources and multiple medical conditions, employing a range of devices coupled with telemedicine and telehealth services. This data collection group continually gathers, organizes, and manages clinical data related to patients. It facilitates categorizing streaming data according to the relevant domain (e.g., specific medical conditions), where records are subsequently published. Apache Kafka, operating as publishes-subscribe messaging system designed for distributed streaming, is a central component in this data management strategy. It is built to be a replicated, distributed, and partitioned service. Health monitoring devices feed real-time data into Kafka via Kafka producers. The fundamental concept that Kafka introduces for a stream of records is termed a "topic." Kafka servers use these topics to store incoming messages from publishers for a defined period before releasing them to the relevant data stream. Each topic is subdivided into multiple partitions, each capable of storing data in diverse formats.

Fig. 3. The workflow of the proposed architecture.



Fig. 4. The Kafka communications system.

Consumers of Kafka access information as it becomes available by subscribing to one or more topics. The Kafka communication infrastructure is depicted in Fig. 4. To ensure the efficient operation of Kafka, ZooKeeper, a centralized service, plays a vital role by providing group services, distributed synchronization, configuration information maintenance, and naming services. Distributed applications use these services extensively, although implementing and maintaining them comes with inherent challenges, such as dealing with recurring defects and race conditions. Typically, applications initially underinvest in these services due to the complexity of their implementation, rendering them fragile in the face of change and difficult to manage. Consequently, resolving these issues and enhancing the robustness of distributed applications is an ongoing endeavor within distributed systems.

This case study involves two data producer programs that simulate connected devices, utilizing Apache Kafka to generate data events. Apache Spark, an open-source, high-speed distributed processing engine, plays a central role in this system. Spark's most notable feature is its capability for in-memory calculations, significantly enhancing its processing speed. Furthermore, Spark offers user-friendly features, an advanced framework for large-scale analysis, and the ability to execute disk-based computing when dealing with datasets that exceed available memory. A key concept employed by Spark is Resilient Distributed Datasets (RDDs), which are distributed, immutable collections of items. To achieve parallelization, Spark internally spreads the RDD data across multiple nodes within the cluster. RDDs can store input and intermediate data in memory, reducing the cost of input-output operations associated with reading from or writing to system files. This

feature enables efficient data reuse, which is particularly beneficial for iterative machine learning algorithms. Once data is transformed into an RDD, two fundamental types of operations can be performed:

- Transformations: These operations involve applying mapping, filtering, and more to existing RDDs to generate new RDDs.

- Actions: These operations compute a result using an RDD, which is then returned or saved to an external storage system.

Spark also includes an ML library, MLlib, which encompasses popular machine learning techniques such as classification, regression, clustering, and more. To handle real-time data from sources like Kafka and Twitter, Spark streaming builds upon the Spark API. The batch-processing Spark engine divides incoming data streams into less than one-second segments, creating discretized streams (DStreams) as high-level abstractions. Each mini-batch within the DStream collection is patterned after a Spark RDD. In this study, Spark is employed for streaming data processing, with Spark streaming managing the Kafka data stream, and MLlib is used to implement machine learning algorithms. Spark adheres to a master-worker architecture for distributed processing. Each Spark application can establish one master process, the executor in Spark, and several worker processes referred to as drivers. These drivers, like the master, are responsible for evaluating, allocating, scheduling, and supervising the tasks among the executors. The driver also maintains the necessary data consistency throughout the program. In contrast, the executors are solely responsible for executing the code assigned by the driver and transmitting the results back to the driver, as depicted in Fig. 5. This architecture ensures the efficient distribution of tasks and data processing within the Spark application.



Fig. 5.   Master-worker architecture.

The classification of data collected from diverse sources for various diseases necessitates using classification models capable of discerning user characteristics in the presence or absence of a disease. In this research, two classification models have been employed, each briefly introduced below:

*1) K-Nearest Neighbor (KNN):* KNN is a versatile supervised learning method that can serve as a classification and regression algorithm. It determines the distance between the test data point and all training data points and selects the K training data points closest to the test data. Based on their distances, the test data point is then assigned to the class that

most of these K neighbors belong to. This method, represented in Method 1, can be briefly described by Algorithm 1.

---

**Algorithm 1 KNN Algorithm**

---

```
1: procedure KNN(Instance, TestData, K)
2:          C ← Size(TestData)
3:          Dist[C][2] ← 0
4:          for i in TestData do
5:              d ← EclideanDistance(i, Instance)
6:              Dist[i][1] ← d
7:              Dist[i][2] ← Class(i)
8:          end for
9:          Srt ← Sort(dist[:][1])  ▷ Sort 2nd column based on that
10:         Sel ← Srt[1 : K][2]
11:         Cls ← Mode(Sel)
12:         return Cls
13: end procedure
```

---

*2) Support Vector Machine (SVM):* SVM is another supervised machine learning method primarily used for classification, although it can also handle regression tasks. In SVM, each data point is represented as a point in an n-dimensional space, where "n" represents the number of features available for classification. Each feature corresponds to a specific coordinate within this space. SVM aims to identify the hyperplane that optimally separates the two classes in the data. Support vectors represent individual data points within this multi-dimensional space, and the SVM classifier seeks to identify the hyperplane or line that maximally divides the two classes. To achieve this, the SVM algorithm considers certain assumptions about the data, aiming to find the best hyperplane:

- Maximizing margin: SVM strives to find the hyperplane that maximizes the margin or the distance between the hyperplane and the nearest data points of both classes. This maximized margin ensures robust separation.

- Support vectors: The data points closest to the hyperplane, known as support vectors, significantly influence the determination of the optimal hyperplane.

- Kernel functions: SVM can employ kernel functions to map the data into higher-dimensional spaces when a linear separation is not feasible. These functions allow SVM to perform non-linear classification effectively. As a classification algorithm, SVM provides the means to efficiently distinguish between different classes within a dataset by defining the most appropriate hyperplane or decision boundary.

The suggested architecture addresses the complex issues involved in forecasting real-time health status in healthcare scenarios powered by the IoT. The applicability of this is emphasized by numerous essential features designed to tackle these particular challenges. The architecture's scalability is a fundamental aspect that allows it to easily handle large and growing datasets often encountered in healthcare. The ability to effortlessly increase resources with data expansion guarantees

consistent performance. The ability to analyze data in real-time is another important aspect, allowing for quick intake, analysis, and understanding of streaming healthcare data.

Furthermore, the architecture's distinctive advantage resides in its implementation of distributed machine learning models specifically created to handle the vastness and complexities of medical data. This enables the simultaneous execution of tasks to enhance the efficiency of training models, hence improving the accuracy and speed of health status forecasts. Moreover, the architecture has exceptional proficiency in incorporating various IoT devices and dissimilar data sources, merging distinct data streams for thorough analysis. By prioritizing security and privacy safeguards, adapting to different data speeds from IoT sensors, and maximizing resource efficiency, it effectively tackles the complex difficulties often seen in healthcare situations powered by IoT. In conclusion, these architectural characteristics together enable the system to effectively negotiate the intricacies of real-time health status prediction, establishing it as an optimal framework for handling the distinct requirements of healthcare data analysis in IoT contexts.

### IV. SCENARIO DESCRIPTIONS

In Scenario 1, Fig. 6 illustrates three hyperplanes labeled A, B, and C. The key principle to selecting the appropriate hyperplane is to choose the one that best separates the two classes. In this scenario, hyperplane B does an excellent job of achieving this separation.



Fig. 6. Three sample hyper-planes.

Scenario 2 presents three hyperplanes (A, B, and C) in Fig. 7. The goal is to choose the hyperplane that maximizes the distance between the closest data point of any class and the hyperplane. This distance is referred to as the margin, as shown in Fig. 8. Hyperplane A has a larger margin than B and C, making it the right choice. Opting for a hyperplane with a larger margin enhances robustness and minimizes the chances of misclassification.

Fig. 7.   Three hyper-planes that could separate two classes.



Fig. 8.   Comparison of three hyper-planes with margins in scenario 2.

In Scenario 3, although hyperplane B has a larger margin than A, SVM prioritizes proper classification of the classes before maximizing the margin. Hyperplane B makes a classification error, whereas A correctly categorizes everything. Therefore, hyperplane A is selected as the appropriate choice (see Fig. 9).



Fig. 9.   Evaluation of hyperplanes A and B in scenario 3.

Scenario 4 involves an outlier, represented by the star, residing in the region of the circle class, making it impossible to separate the two classes using a straight line. However, the SVM algorithm can disregard outliers and identify the

hyperplane with the maximum margin. As a result, SVM classification is robust against outliers (see Fig. 10).



Fig. 10.  Robustness of SVM against outliers in scenario 4.

In Scenario 5, when a linear hyperplane is insufficient to categorize two classes, a new feature, $z = x^2 + y^2$, is introduced to create a three-dimensional representation of the data points, as shown in Fig. 11 and Fig. 12. This new feature, z, is a mathematical construct that enables the creation of a linear hyperplane, making it possible for SVM to classify the two classes effectively. The SVM algorithm employs the "kernel trick" to automatically find this hyperplane. The SVM kernel is a function that transforms non-separable problems into separable ones by projecting data from a low-dimensional input space into a higher-dimensional space. This is particularly valuable for addressing problems with non-linear separations. It performs intricate data transformations before determining how to split the data based on the provided labels or outputs. The hyperplane appears as a circle in the original input space, as depicted in Fig. 13. The kernel trick allows SVM to handle complex, non-linear separations and enables the classification of data that cannot be linearly separated in the original feature space.



Fig. 11.  Introduction of a new feature in scenario 5.

Fig. 12. Three-dimensional representation with additional feature (z) in scenario 5.



Fig. 13. Application of SVM kernel trick for non-linear separation in scenario 5.

## V. DISCUSSION

To ensure high data availability and avoid a single point of failure, it is essential to store the results and data streams generated by each user in a distributed manner. Distributed databases outperform traditional database systems in terms of performance and scalability. Apache Cassandra is an open-source, distributed, and free NoSQL database system designed to handle massive volumes of data, whether structured, unstructured, or semi-structured, across multiple computers. Cassandra's architecture greatly enhances its scalability, operational capabilities, and continuous accessibility. It also offers rapid write and read rates when used with Spark. Distributed databases provide several valuable features:

- Affordability and ease of use: Distributed databases are cost-effective and straightforward.

- Data transfer speed: They offer significantly faster data transport than traditional databases.

- Scalability: Distributed databases can be scaled easily by adding columns, accelerating the processing of larger and more data.

- Cluster scalability: Distributed databases can expand their cluster capacity by adding more nodes without a specific distribution. After processing data with Spark, the output data is stored in a table using Cassandra and a primary key. This database can be accessed later for real-time monitoring, reporting, and analysis of historical data.

- Data replication and partitioning: Data is replicated across various computers to enhance data availability and fault tolerance.

TABLE I.        UCI HEART DISEASE DATASET

| No | Attribute No | Attribute Name | Description |
|---|---|---|---|
| 1 | 3 | Age | Age of Patients |
| 2 | 4 | Sex | 0/1(M/F) |
| 3 | 9 | CP | Type of Chest Pain |
| 4 | 10 | TRestBPS | Blood Pressure when the Patient is on |
|  |  |  | Rest |
| 5 | 12 | Chol | Blood Cholesterol |
| 6 | 16 | FBS | Fasting Blood Sugar |
| 7 | 19 | RestECG | ElectroCardioGraphic when Patient |
|  |  |  | is on Rest |
| 8 | 32 | Thalach | Heart Rate(Max) |
| 9 | 38 | Exang | Exercise Causes Angina (Y/N = |
|  |  |  | 1/0) |
| 10 | 40 | OldPeak | Exercise-induced ST depression in |
|  |  |  | comparison to rest |
| 11 | 41 | Slope | The Peak Exercises in cline ST section. (UpSloping/Flat/DownSloping = 1/2/3 |
| 12 | 44 | CA | Main Vessels Colored with |
|  |  |  | Fluoroscopy in Number (0–3) |
| 13 | 51 | Thal | Normal/ Fixed defect/ Reversible |
|  |  |  | Defect = 3/6/7 |
| 14 | 58 | Num(Class) | Heart Disease Diagnosis (Status of |
|  |  |  | Angiographic Disease) if Diameter |
|  |  |  | Narrowing¡= 50% =0 Otherwise =1 |

Apache Zeppelin is an open-source data analysis environment that works with Apache Spark. It is a web-based, versatile notebook that facilitates interactive data analysis, real-time data exploration, visualization, and collaboration. Zeppelin supports an expanding list of programming languages and interfaces, including SparkSQL, Hive, AngularJS, Scala, Python, markdown, and Shell. Using Scala, it can create

dynamic, data-driven, and collaborative documents, among other capabilities. Apache Zeppelin is valuable for writing, organizing, and executing analytical code and visualizing results across extensive workflows. Zeppelin can automatically generate input forms in your notebook, provide simple visualizations to present results, and allow colleagues to share the notebook's URL. In real-time data retrieval from the Cassandra database, a Zeppelin dashboard is developed to display data in charts, tables, and other formats. This dashboard updates its data every second, allowing authorized individuals, such as doctors, healthcare companies, or external consultants, to access the data regardless of their patient's or client's health status.

In this research, two datasets obtained from well-known data sources, Kaggle and UCI, were utilized. These datasets pertain to medical conditions, specifically diabetes and heart diseases. Table I provides an overview of the information related to these datasets. It is worth noting that although the Cleveland dataset contains 76 attributes, previous studies have primarily focused on using only a subset of fourteen attributes. Among the various datasets, the Cleveland dataset has been the primary focus of machine learning researchers. The "Class" field in these datasets indicates the presence or absence of a particular medical condition, such as heart disease. The values in the "Class" field range from zero (indicating no presence of the condition) to four, with the Cleveland dataset primarily concentrating on discriminating between the presence (values 1, 2, 3, 4) and absence (value 0) of heart disease.

The dataset used in this research was sourced from Kaggle's Diabetes Dataset. Kaggle is a well-known platform for data science competitions and provides a freely available dataset that numerous authors have used in previous studies. This dataset consists of ten features and 15,000 observations, and it is employed to predict whether a patient has diabetes. Table II offers an overview of the features included in this dataset.

TABLE II.    KAGGLE DIABETES DATASET DESCRIPTIONS

| No | Attribute Name | Description |
|---|---|---|
| 1 | Patient ID | Patient Identification Number |
| 2 | Pregnancies | A patient gets diabetes after Pregnancy |
| 3 | Plasma Glucose | Glucose amount in Blood |
| 4 | Diastolic Blood Pressure | Blood Pressure when Patient is on Rest |
| 5 | Triceps Thickness | Body Fat |
| 6 7 8 | Serum Insulin Body Mass Index(BMI) Diabetes Pedigree | Insulin amount in Blood $W\ eightinKG\ Height2\ inM\ 2$ ) Diabetes History in Family |
| 9 | Age | Patient Age |
| 10 | Class | Diabetic = 1, NonDibaetic = 0 |

## VI.    RESULTS

The proposed real-time health status forecasting system is driven by a single-node cluster featuring a Core i7 CPU, 16 GB of RAM, and the Ubuntu 20.04 operating system. This system seamlessly integrates the trained model with Kafka streaming data processing and runs on the Spark platform. As depicted in Fig. 14, the application establishes a connection to Kafka streaming and commences receiving data streams from various Kafka producers. When it encounters streams related to health characteristics, it retrieves the attribute values from each topic within the illness events sent via Kafka streaming. Subsequently, it employs the trained model to predict the health state of the individuals. In parallel, the Cassandra database records each forecasted health state in a table, employing the identification (ID) as the primary key, which is ideal for ensuring data redundancy and reliability. This stored data can later be queried to examine historical information.



Fig. 14.  Apply classification algorithms.

All the tests were conducted using a cluster configuration consisting of one primary node and two worker nodes, each running the Ubuntu 20.4 operating system within VMware virtual environments. Several steps were undertaken to

facilitate communication between the nodes and ensure the proper functioning of the Spark application:

- User accounts and development environment setup: Spark user accounts were created to simplify inter-node communication. Scala and Java were installed. Open SSH Server was set up. Key pairs were generated to enable passwordless SSH configuration across the nodes, ensuring that the Spark master can effectively connect, launch, pause, and run tasks on multiple worker nodes.

- Software installation: Spark, Kafka, and Cassandra were unpacked and installed on a single node. Two themes and corresponding tables were created, one for diabetic disease and the other for heart illness.

- Environment variables: The bashrc file was modified to include essential environment variables like SPARK and JAVA_HOME in the home directory.

- Node replication: To ensure uniformity and consistency across various nodes, the setup folder of the single-node cluster was duplicated multiple times, with one node designated as the master and the others as workers.

- Hostname and host configuration: The hostname and hosts were modified on all nodes to facilitate proper inter-node communication.

The primary stages for implementing the Spark application in a Zeppelin notebook are as follows:

- Spark context and streaming context creation: An instance of Spark contexts and streaming context was created to access all Spark streaming functionalities.

- Direct stream creation: A direct stream was created using the specified Kafka parameters and topics.

- Data extraction: The identifiers and characteristics of each topic and stream were extracted.

- ML model utilization: The pre-trained ML model was used to predict the health status.

- Data storage: All attributes and the predicted labels were saved to the Cassandra keyspace and table.

- Streaming start: The Spark streaming context was initiated using the start method, allowing real-time health data processing.

The research allocated 25% of the data for testing purposes, while the remaining 75% was used to train the machine learning models. The datasets were divided into training and test datasets randomly. To address the issue of the computational cost of sorting feature values across large distributed datasets, an approximate set of candidate splits was identified over a sampled portion of the data. This method has

been shown to enable more accurate predictions by analyzing the model error and the test data, effectively mitigating the negative impacts of both underfitting and overfitting. One of the most crucial and valuable measures for assessing the performance of testing and treatment is the Receiver Operating Characteristic (ROC) curve. The MLlib provides support for ROC curve evaluation. On the other hand, classification accuracy is determined by the ratio of all correct predictions to all the prediction data. The classification accuracy for the datasets in this study was assessed using the following equation:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \tag{1}$$

Sensitivity and specificity are two critical metrics used to assess classification models' performance, particularly in medical diagnoses and other fields where accurate predictions are crucial. These metrics are calculated as follows:

Sensitivity (True Positive Rate or Recall): Sensitivity is the percentage of actual positives (e.g., patients with a specific condition) that are correctly identified by the model. It indicates the model's ability to detect true positive cases.

Specificity: Specificity is the percentage of actual negatives (e.g., patients without the condition) that are correctly identified by the model. It measures the model's ability to avoid false positive predictions.

In these equations, true positives (TPs) are the cases correctly classified as positive, true negatives (TNs) are the cases correctly classified as negative, false positives (FPs) are cases incorrectly classified as positive (when they are actually negative), and false negatives (FNs) are cases incorrectly classified as negative (when they are actually positive).

Sensitivity and specificity provide insights into the model's performance in terms of both correctly identifying individuals with the condition and correctly identifying individuals without the condition. Balancing these two measures is important, especially in situations, where missing a true positive (e.g., a medical condition) or incorrectly identifying a false positive (unnecessary treatment or diagnosis) has significant consequences.

$$Sensitivity = \frac{TP}{TP+FN} \tag{2}$$

$$Specifity = \frac{TN}{TN+FP} \tag{3}$$

Our machine learning model's effectiveness was assessed on two established datasets. The empirical results indicate that our utilization of Spark for the execution of the proposed methodology demonstrates notable efficiency and scalability (see Fig. 14 and Fig. 15). Fig. 16 shows the specificity, sensitivity, ROC curve and accuracy obtained in heart disease. The findings further underscore that the proposed model consistently delivers dependable and superior predictive outcomes.

Fig. 15. Specificity, sensitivity, ROC Curve, and accuracy obtained in diabetes datasets in comparison to the other algorithms.



Fig. 16. Specificity, sensitivity, ROC Curve, and accuracy obtained in heart disease datasets in comparison to other algorithms.

## VII. CONCLUSION

This study has demonstrated the successful application of a machine learning model for real-time health status prediction in the healthcare domain. By employing Apache Spark in conjunction with Kafka streaming and Cassandra, we have created an efficient and scalable system for processing and analyzing healthcare data streams. The results of our empirical tests on two distinct datasets reveal that our model consistently provides reliable and high-quality predictions. However, the current design also reveals inherent constraints in the scalability of traditional data storage systems like Cassandra when handling exponentially growing healthcare data. While effective for many use cases, these systems might face challenges in handling future data volume surges, potentially leading to performance bottlenecks and increased resource requirements. The ability to monitor and predict health conditions in real-time is of paramount importance, particularly in the context of chronic illnesses and emergencies. Our proposed system offers a promising solution for continuous health monitoring and early detection, potentially saving lives and reducing healthcare costs. The key takeaway from our research is the effectiveness of combining advanced technologies like Spark, Kafka, and Cassandra to process and analyze healthcare data streams. This approach opens up new possibilities for healthcare analytics and real-time monitoring, benefiting patients, healthcare providers, and the broader medical community. In the future, we envision further refinements and enhancements to our system, including the integration of additional data sources and the development of more sophisticated machine-learning algorithms. As the healthcare sector continues to generate vast amounts of data, the need for innovative solutions like the one presented in this study will only grow, ushering in a new era of data-driven healthcare.

## REFERENCES

[1] M. Asch et al., "Big data and extreme-scale computing: Pathways to convergence-toward a shaping strategy for a future software and data ecosystem for scientific inquiry," The International Journal of High Performance Computing Applications, vol. 32, no. 4, pp. 435-479, 2018.

[2] M. El Samad, S. El Nemar, G. Sakka, and H. El-Chaarani, "An innovative big data framework for exploring the impact on decision-making in the European Mediterranean healthcare sector," EuroMed Journal of Business, vol. 17, no. 3, pp. 312-332, 2022.

[3] M. Karatas, L. Eriskin, M. Deveci, D. Pamucar, and H. Garg, "Big Data for Healthcare Industry 4.0: Applications, challenges and future

perspectives," Expert Systems with Applications, vol. 200, p. 116912, 2022.

[4] S. S. Ghahfarrokhi and H. Khodadadi, "Human brain tumor diagnosis using the combination of the complexity measures and texture features through magnetic resonance image," Biomedical Signal Processing and Control, vol. 61, p. 102025, 2020.

[5] K. Batko and A. Ślęzak, "The use of Big Data Analytics in healthcare," Journal of big Data, vol. 9, no. 1, p. 3, 2022.

[6] A. A. Zein, S. Dowaji, and M. I. Al-Khayatt, "Clustering-based method for big spatial data partitioning," Measurement: Sensors, vol. 27, p. 100731, 2023.

[7] M. Mohtasebi et al., "Detection of low-frequency oscillations in neonatal piglets with speckle contrast diffuse correlation tomography," Journal of Biomedical Optics, vol. 28, no. 12, pp. 121204-121204, 2023.

[8] R. Nathan et al., "Big-data approaches lead to an increased understanding of the ecology of animal movement," Science, vol. 375, no. 6582, p. eabg1780, 2022.

[9] C. Acciarini, F. Cappa, P. Boccardelli, and R. Oriani, "How can organizations leverage big data to innovate their business models? A systematic literature review," Technovation, vol. 123, p. 102713, 2023.

[10] M. Andronie et al., "Big Data Management Algorithms, Deep Learning-Based Object Detection Technologies, and Geospatial Simulation and Sensor Fusion Tools in the Internet of Robotic Things," ISPRS International Journal of Geo-Information, vol. 12, no. 2, p. 35, 2023.

[11] W. Li, "Big Data precision marketing approach under IoT cloud platform information mining," Computational intelligence and neuroscience, vol. 2022, 2022.

[12] M. Q. Bashabsheh, L. Abualigah, and M. Alshinwan, "Big data analysis using hybrid meta-heuristic optimization algorithm and MapReduce framework," in Integrating meta-heuristics and machine learning for real-world optimization problems: Springer, 2022, pp. 181-223.

[13] G. S. Bhathal and A. Singh, "Big Data: Hadoop framework vulnerabilities, security issues and attacks," Array, vol. 1, p. 100002, 2019.

[14] A. Adnan, Z. Tahir, and M. A. Asis, "Performance evaluation of single board computer for hadoop distributed file system (hdfs)," in 2019 International Conference on Information and Communications Technology (ICOIACT), 2019: IEEE, pp. 624-627.

[15] S. Vairachilai, A. Bostani, A. Mehbodniya, J. L. Webber, O. Hemakesavulu, and P. Vijayakumar, "Body Sensor 5 G Networks Utilising Deep Learning Architectures for Emotion Detection Based On EEG Signal Processing," Optik, p. 170469, 2022.

[16] M. Bolhassani and I. Oksuz, "Semi-Supervised Segmentation of Multi-vendor and Multi-center Cardiac MRI," in 2021 29th Signal Processing and Communications Applications Conference (SIU), 2021: IEEE, pp. 1-4.

[17] S. P. Rajput et al., "Using machine learning architecture to optimize and model the treatment process for saline water level analysis," Journal of Water Reuse and Desalination, 2022.

[18] S. R. Abdul Samad et al., "Analysis of the Performance Impact of Fine-Tuned Machine Learning Model for Phishing URL Detection," Electronics, vol. 12, no. 7, p. 1642, 2023.

[19] V. Monjezi, A. Trivedi, G. Tan, and S. Tizpaz-Niari, "Information-Theoretic Testing and Debugging of Fairness Defects in Deep Neural Networks," arXiv preprint arXiv:2304.04199, pp. 1571-1582, 2023 2023, doi: 10.1109/ICSE48619.2023.00136.

[20] W. Anupong et al., "Deep learning algorithms were used to generate photovoltaic renewable energy in saline water analysis via an oxidation process," Water Reuse, vol. 13, no. 1, pp. 68-81, 2023.

[21] M. Khodayari, J. Razmi, and R. Babazadeh, "An integrated fuzzy analytical network process for prioritisation of new technology-based firms in Iran," International Journal of Industrial and Systems Engineering, vol. 32, no. 4, pp. 424-442, 2019.

[22] J. Zandi, A. N. Afooshteh, and M. Ghassemian, "Implementation and analysis of a novel low power and portable energy measurement tool for wireless sensor nodes," in Electrical Engineering (ICEE), Iranian Conference on, 2018: IEEE, pp. 1517-1522, doi: 10.1109/ICEE.2018.8472439.

[23] Y. Lu, Z. Miao, P. Sahraeian, and B. Balasundaram, "On atomic cliques in temporal graphs," Optimization Letters, vol. 17, no. 4, pp. 813-828, 2023.

[24] M. R. Moradi, S. R. N. Kalhori, M. G. Saeedi, M. R. Zarkesh, A. Habibelahi, and A. H. Panahi, "Designing a Remote Closed-Loop Automatic Oxygen Control in Preterm Infants," Iranian Journal of Pediatrics, vol. 30, no. 4, 2020.

# A Yolo-based Approach for Fire and Smoke Detection in IoT Surveillance Systems

Dawei Zhang*

College of Information Engineering, Liaodong University, Dandong 118000, Liaoning, China

*Abstract*—**Fire and smoke detection in IoT surveillance systems is of utmost importance for ensuring public safety and preventing property damage. While traditional methods have been used for fire detection, deep learning-based approaches have gained significant attention due to their ability to learn complex patterns and achieve high accuracy. This paper addresses the current research challenge of achieving high accuracy rates with deep learning-based fire detection methods while keeping computation costs low. This paper proposes a method based on the Yolov8 algorithm that effectively tackles this challenge through model generation using a custom dataset and the model's training, validation, and testing. The model's efficacy is succinctly assessed by the precision, recall and F1-curve metrics, with notable proficiency in fire detection, crucial for early warnings and prevention. Experimental results and performance evaluations show that our proposed method outperforms other state-of-the-art methods. This makes it a promising fire and smoke detection approach in IoT surveillance systems.**

*Keywords—IoT; surveillance systems; fire detection; deep learning; Yolov8*

## I. INTRODUCTION

The capacity of Internet of Things (IoT) surveillance systems to track and analyze data in real-time for a range of applications has attracted a lot of interest in recent years [1-3]. One such use is the detection of fire and smoke in hospitals, where quick action is essential to avoid property damage and human casualties [4, 5].

Sensor-based systems, image processing, and machine learning are now used in IoT monitoring systems for fire and smoke detection [6]. Among these innovations, computer vision-based systems have demonstrated encouraging results and drawn the attention of several researchers due to their capacity to deliver precise and trustworthy detection findings [7, 8]. Many methods for improving the performance of these systems have been the subject of recent investigations. Despite advancements in computer vision-based fire and smoke detection systems, there are still a number of restrictions and research gaps [9, 10]. False alarms, poor precision, and a restricted capacity to adjust to changing circumstances are some of these drawbacks [11]. These issues indicate the need for more study to enhance the effectiveness and performance of these systems.

In IoT surveillance systems, recent developments in deep learning-based techniques have shown promise in terms of the precision and accuracy [11-14]. Several studies have suggested deep learning-based techniques to solve the shortcomings of current computer vision-based systems, particularly employing the YOLO algorithm. These experiments have shown considerable gains in detection speed and accuracy over conventional machine learning methods.

The advancement of deep learning-based approaches for fire and smoke detection is impeded by several challenges, notably the scarcity of adequate training data and the intricate nature of environmental variables [15, 16]. Overcoming these obstacles is pivotal to enhancing the reliability and effectiveness of deep learning systems dedicated to fire and smoke detection. A critical avenue for improvement lies in conducting in-depth research and comprehensive analysis to delve into the intricacies of these challenges [17]. By addressing the issues related to insufficient training data and navigating the complexities of diverse environmental factors, researchers can refine and optimize deep learning models. This iterative process is essential for fortifying the robustness of these systems, ultimately paving the way for more accurate and dependable fire and smoke detection in various real-world scenarios [18]. As the field progresses, a concerted effort to explore and resolve these challenges will contribute significantly to the evolution of deep learning methodologies, fostering advancements that hold substantial promise for applications in safety and security within IoT surveillance systems.

To address these challenges, this study proposes a deep learning-based approach using the YOLOv8 algorithm for fire and smoke detection in health houses. Using a standard dataset for training, validation, and testing processes, which includes a diverse range of scenarios and environments. The proposed method is evaluated on the custom dataset and compared to existing state-of-the-art methods, demonstrating superior performance in terms of accuracy and speed. The main contributions of this research are as follows:

*1) Identifying* the limitations and research gaps in existing computer vision-based fire and smoke detection systems in health houses.

*2) Proposing* a deep learning-based approach using the YOLOv8 algorithm to address these challenges as well as improve the accuracy and speed of detection.

*3) Evaluating* the proposed method on a custom dataset using extensive performance evaluation metrics.

The reminder of this paper is as, Section II review of previous studies. Section III discuss about material and methods. Section IV presents results and discussions. Finally, this paper concludes in Section V.

## II. RELATED WORKS

Convolutional neural networks (CNNs) were proposed in the study [19] as a technique for early fire detection during surveillance for efficient disaster management. The method involves acquiring surveillance footage, pre-processing the footage, and training a CNN to detect fires. The proposed method achieves an accuracy of 97.4% for fire detection, which is higher than existing methods. Key features of the method include the use of CNNs for accurate fire detection and the ability to detect fires in real-time. The suggested approach can be used in various scenarios, including forest fires, building fires, and wildfires. Limitations of the method include the need for high-quality surveillance footage and the potential for false positives in certain scenarios.

The paper in [20] presented a deep learning-based forest fire detection approach utilizing unmanned aerial vehicles (UAVs) and the YOLOv3 object detection algorithm. The method consists of acquiring high-resolution images from UAVs, pre-processing the images, detecting the presence of fire using YOLOv3, and sending the location of the fire to a control center. The proposed method achieves an accuracy of 96.4% for forest fire detection. Key features of the method include the use of UAVs for acquiring high-resolution images and the use of YOLOv3 for object detection. The proposed approach has the potential for use in real-world scenarios, but its limitation is the dependency on good weather conditions for successful operation.

The authors in [21] present a deep learning framework called Fire-Net for active forest fire detection. The method involves acquiring high-resolution images from a network of ground-based cameras, pre-processing the images, and training a deep neural network to detect fires. Fire-Net achieves an accuracy of 98.7% for fire detection, which is higher than existing methods. Key features of the method include the use of high-resolution images and a deep neural network for accurate fire detection. The proposed approach has the potential for use in real-world scenarios, but its limitation is the dependency on good weather conditions for successful operation.

The authors in [22] proposed an improved forest fire detection method using a deep learning approach and the Detectron2 model. The method involves acquiring high-resolution images, pre-processing the images, training the Detectron2 model, and detecting fires. The proposed approach achieves an accuracy of 98.6% for fire detection, which is higher than existing methods. Key features of the method include the use of the Detectron2 model for accurate fire detection and the ability to detect fires in real-time. The proposed approach has the potential for use in various scenarios, including forest fires, building fires, and wildfires. Limitations of the method include the need for high-quality images and the potential for false positives in certain scenarios.

The paper in [12] presented a forest fire notification and detection method using IoT and AI approaches. The method involves installing IoT sensors in forested areas to detect environmental conditions and using a deep learning algorithm to detect fires. Key features of the method include the use of IoT sensors for environmental monitoring and a deep learning algorithm for accurate fire detection. The proposed approach has the potential for use in real-world scenarios, but its limitation is the high cost associated with installing and maintaining IoT sensors.

The authors in [23] proposed a wildfire and smoke detection method using ensemble CNN and a staged YOLO model. The method involves acquiring high-resolution images from a network of cameras, pre-processing the images, and using a staged YOLO model for smoke detection and an ensemble CNN for wildfire detection. Key features of the method include the use of a staged YOLO model and ensemble CNN for accurate detection of smoke and wildfire, respectively. The proposed approach has the potential for use in real-world scenarios, but its limitation is the need for high-quality images and the potential for false positives in certain scenarios.

## III. MATERIAL AND METHOD

This section provides an overview of the study's content and methodology. Yolov8 is the foundation of the primary methodology employed in this study. Initially, YOLO was a pre-trained object detector programmed to recognize commonplace items such as chairs, tables, phones, cars, etc. This study presents a detection technique according to YOLO algorithms to create a model that might identify. In real-time applications, the models perform well as well.

### A. Dataset

The dataset consists of three different types of jewelry (background, fire, and smoke). The dataset contains pictures captured by webcam. The obtained images from a webcam permanently installed at a jewelry store. Some samples of the dataset are shown in Fig. 1. Small target items, which are harder to detect, are included in the collection, along with images of various sizes. The data was collected from GitHub resource[1].

In order to prepare the dataset to be more robust, augmenting of the images in the dataset is performed. Choosing images with various angles, sizes, resolutions, forms, and sample counts in each image. A maximum of three augmented versions of each image were created by applying the following random effects: horizontal flip, rotation between -15° and +15°, exposure between -10% and +10%, brightness between -20% and +20%, and saturation between -20% and +20%. Photos that have been enhanced. The final step is to separate the labeled images into a validation set (6%) and a test set (1%), training set (93%).

### B. Google Colab

Using Google Colab, which offers free usage of potent GPUs. All testing and training tasks are performed utilizing a 12GB NVIDIA Tesla T4 GPU; more details are given in Fig. 2. All the models are trained for 50 epochs with an image size of 640 and with YOLO default adjustment for other hyperparameters.

---

[1]https://github.com/Abonia1/YOLOv8-Fire-and-Smoke-Detection/tree/main/ datasets/fire-8

Fig. 1.    Sample images from the dataset.



Fig. 2.    Details of Google colab's GPU.

*C. Yolov8 Model*

The latest version of YOLO is Yolov8, which was released in 2022 by Ultralytics [24]. YOLOv8 is an object detection model that is based on the You Only Look Once (YOLO) family of algorithms. The architecture of Yolov8 is shown in Fig. 3. YOLOv8 has several architectural improvements over the previous versions, such as:

*1) Convolutional neural network (CNN) architecture:* Yolov8 uses a modified version of the EfficientNet backbone network, which is a cutting-edge CNN design that strikes a solid balance between computing efficiency and accuracy. Yolov8's backbone network is in charge of taking features out of the input image. Yolov8's backbone network is specifically based on a modified version of the EfficientNet design. Modern convolutional neural network (CNN) architecture EfficientNet aims to strike a reasonable compromise between accuracy and processing efficiency. A succession of convolutional layers that successively reduce the spatial resolution of the feature maps while increasing the number of channels make up Yolov8's EfficientNet backbone. A sizable image classification dataset, such as ImageNet, is used to pre-train the backbone network.

*2) Feature aggregation:* Yolov8 uses a feature pyramid network (FPN) to aggregate features at different scales, which improves the model's ability to detect objects of different sizes.

*3) Object detection head:* The head is made up of some convolutional layers, followed by an output layer that generates the results of the detection. Each object in the input image is predicted to have bounding boxes, objectness scores, and class probabilities by Yolov8's object detection head. The head is made up of some convolutional layers, followed by an output layer that generates the results of the detection. Yolov8's head is made up of three prediction branches that, in turn, forecast bounding boxes, objectness scores, and class probabilities. The implementation of each branch consists of a set of convolutional layers, followed by an output layer that generates the corresponding predictions. The number of anchor boxes and object classes determines the number of output channels in each branch.

*4) Training strategy:* Yolov8 enhances the accuracy and speed of the model by combining anchor-based and anchor-free item detection techniques. Additionally, a progressive scaling strategy is used during training to enhance the model's capacity to recognize objects of various sizes.

Fig. 3.   Yolov8 architecture.

## IV. RESULTS AND DISCUSSION

### A. Experimental Results

To conduct experimental results for fire and smoke detection using a YOLOv8 model, sample output images from both validation and testing sets can be examined. These images can be obtained by running the trained model on the validation and testing datasets and visualizing the resulting predictions. For each image, the model correctly detects whether there is a fire or smoke present. Fig. 4 shows the samples of experimental results.

### B. Performance Measurements

This study uses precision, recall, and F1socre criteria to assess performance. The P-curve, R-curve, PR-curve, and F1-curve are assessment metrics frequently employed in machine learning to gauge the effectiveness of models [16, 18]. The fraction of true positive predictions among all positive predictions is represented by the P-curve in the case of a fire and smoke detection Yolov8 model. The R-curve shows the recall of the model or the percentage of accurate positive predictions out of all actual positive samples. The link between accuracy and recall is depicted by the PR-curve, which also illustrates how effectively the model balances the two parameters. Finally, the F1-curve provides a comprehensive evaluation of the model's performance by combining precision and recall into a single score.

Fig. 4.    Samples of experimental results.



Fig. 5.    P-curve of the model.

As shown in Fig. 5, the graph displays the YOLOv8 model's validation set's fire and smoke detection precision values. The precision values are represented on the y-axis, while various confidence criteria, ranging from 0.0 to 1.0, are represented on the x-axis. The fraction of accurate positive forecasts among all positive predictions is known as the precision value. The graph demonstrates that the accuracy of detecting fire is always greater than the accuracy of detecting smoke. The precision value for detecting a fire is approximately 0.8 at a confidence threshold of 0.5, whereas the precision value for detecting smoke is approximately 0.5. These results indicate that the model performs better at identifying fire than smoke.

As shown in Fig. 6, the recall values of a YOLOv8 model for fire and smoke detection on the validation data are shown in the R-curve graph. The recall values are displayed on the y-axis, and the various confidence criteria are displayed on the x-axis. The recall is the percentage of correctly predicted positive samples among all truly positive samples. According to the graph, the model can detect fires more often than smoke, with a maximum recall of 0.9 for fire and 0.7 for smoke at a confidence level of 0.5. This implies that the model is more effective at identifying fire than smoke, while changes could also influence this in the visual properties of the two classes or the training data.



Fig. 6.   R-curve of the model.



Fig. 7.   PR-curve of the model.

Fig. 7 depicts the precision-recall trade-off for a YOLOv8 fire and smoke detection model. The recall values are on the x-axis, and the precision values are on the y-axis. From 0.0 to 1.0, the precision and recall values are measured; higher values denote greater performance. The graph demonstrates that the model generates greater precision values for the identification of fire at all recall levels, demonstrating that it is more accurate in identifying fire than smoke. Additionally, the precision and recall values for smoke detection are considerably lower than those for fire detection, as the curve shows, suggesting that the model may have some difficulties identifying smoke. The PR-curve offers an extensive picture of the model's performance in terms of precision and recall.

As shown in Fig. 8, the harmonic means of recall and precision for fire and smoke detection in the validation set of a YOLOv8 model is displayed on the F1-curve graph. The y-axis displays the F1 score values ranging from 0.0 to 1.0, while the x-axis displays various confidence criteria. The F1 score,

which considers precision and recall, is a frequently used metric for assessing a model's overall performance. The greatest F1 score for fire detection at a confidence level of 0.5 is 0.86, while the maximum F1 score for smoke detection is 0.57. This shows that the model performs better at detecting fire than smoke. The graph demonstrates that at all confidence criteria, the F1 score for fire detection is consistently higher than the F1 score for smoke detection. This can be caused by variances in the two classes' visual qualities or variations in the training data. The model can still identify smoke properly to a certain extent, as seen by the F1 score for smoke detection, which is still rather high.

The model's performance in terms of precision and recall for fire and smoke detection is thus usefully summarized by the F1-curve. The model appears effective at detecting fires, essential for early warning and prevention, as indicated by the high F1 scores for fire detection. However, the model's performance in detecting smoke might still be improved.



Fig. 8. F1-score of the model.

## V. CONCLUSION

In conclusion, this research paper emphasizes the significance of developing accurate and efficient fire and smoke detection systems in IoT surveillance systems for health houses. While traditional methods have limitations, computer vision-based systems have shown promising results, particularly deep learning-based methods using the YOLO algorithm. However, challenges still exist, such as limited training data and complex environmental factors, which require further investigation. The proposed method in this paper addresses these challenges and demonstrates superior performance in terms of accuracy and speed compared to existing state-of-the-art methods. This research contributes to the ongoing efforts to improve the reliability and effectiveness of fire and smoke detection systems in IoT surveillance systems. For future study, while the YOLO algorithm has shown promise in improving the accuracy and speed of fire and

smoke detection, other deep learning algorithms may also provide superior results. Future studies could explore using other deep learning algorithms, such as Faster R-CNN or Mask R-CNN, and compare their performance to the YOLO algorithm. Moreover, future studies could investigate the development of more advanced vision-based sensors, such as multi-spectral or hyperspectral sensors, which can capture a wider range of data and improve the accuracy of fire and smoke detection.

## REFERENCES

[1] Casola, A. De Benedictis, A. Riccio, D. Rivera, W. Mallouli, and E. M. de Oca, "A security monitoring system for internet of things," Internet of Things, vol. 7, p. 100080, 2019.

[2] K. K. Patel, S. M. Patel, and P. Scholar, "Internet of things-IOT: definition, characteristics, architecture, enabling technologies, application & future challenges," International journal of engineering science and computing, vol. 6, no. 5, 2016.

[3] A. A. Mei Choo Ang, Kok Weng Ng, Elankovan Sundararajan, Marzieh Mogharrebi, Teck Loon Lim, "Multi-core Frameworks Investigation on A Real-Time Object Tracking Application," Journal of Theoretical & Applied Information Technology, 2014.

[4] S. Vijayalakshmi and S. Muruganand, "Internet of Things technology for fire monitoring system," Int. Res. J. Eng. Technol, vol. 4, no. 6, pp. 2140-2147, 2017.

[5] S. Vijayalakshmi and S. Muruganand, "A survey of Internet of Things in fire detection and fire industries," in 2017 International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud)(I-SMAC), 2017: IEEE, pp. 703-707.

[6] Z. Deng, Y. Cao, X. Zhou, Y. Yi, Y. Jiang, and I. You, "Toward efficient image recognition in sensor-based IoT: a weight initialization optimizing method for CNN Based on RGB influence proportion," Sensors, vol. 20, no. 10, p. 2866, 2020.

[7] A. A. Mohammed, N. B. Alapaka, C. Gudivada, K. Bharath, and M. Rajesh Kumar, "Computer Vision Based Autonomous Fire Detection and IoT Based Fire Response System," in Proceedings of International Conference on Communication and Computational Technologies: ICCCT 2021, 2021: Springer, pp. 551-560.

[8] H. Yar, T. Hussain, Z. A. Khan, D. Koundal, M. Y. Lee, and S. W. Baik, "Vision sensor-based real-time fire detection in resource-constrained IoT environments," Computational intelligence and neuroscience, vol. 2021, 2021.

[9] F. Bu and M. S. Gharajeh, "Intelligent and vision-based fire detection systems: A survey," Image and vision computing, vol. 91, p. 103803, 2019.

[10] S. Das, J. Das, O. Krishna, and J. Maiti, "Image Processing-Based Fire Detection Using IoT Devices," in Machine Vision for Industry 4.0: CRC Press, 2022, pp. 207-224.

[11] C. J. Ezeofor and N. O. Nwazor, "An IoT-Based Fire Image Recognition for Home/Industry Security Using Machine Learning," International Journal of Advances in Engineering and Management (IJAEM), vol. Volume 4, pp. 834-842, 2022.

[12] K. Avazov, A. E. Hyun, A. A. Sami S, A. Khaitov, A. B. Abdusalomov, and Y. I. Cho, "Forest Fire Detection and Notification Method Based on AI and IoT Approaches," Future Internet, vol. 15, no. 2, p. 61, 2023.

[13] K. Deepa, A. Chaitra, K. Jhansi, A. K. RD, and M. M. Kodabagi, "Development of Fire Detection surveillance using machine learning & IoT," in 2022 IEEE 2nd Mysore Sub Section International Conference (MysuruCon), 2022: IEEE, pp. 1-6.

[14] A. Aghamohammadi et al., "A deep learning model for ergonomics risk assessment and sports and health monitoring in self-occluded images," Signal, Image and Video Processing, pp. 1-13, 2023.

[15] S. Chaturvedi, P. Khanna, and A. Ojha, "A survey on vision-based outdoor smoke detection techniques for environmental safety," ISPRS Journal of Photogrammetry and Remote Sensing, vol. 185, pp. 158-187, 2022.

[16] R. K Mohammed, "A real-time forest fire and smoke detection system using deep learning," International Journal of Nonlinear Analysis and Applications, vol. 13, no. 1, pp. 2053-2063, 2022.

[17] D.-K. Kwak and J.-K. Ryu, "A Study on Fire Detection Using Deep Learning and Image Filtering Based on Characteristics of Flame and Smoke," Journal of Electrical Engineering & Technology, pp. 1-9, 2023.

[18] L. Hu, C. Lu, X. Li, Y. Zhu, Y. Lu, and S. Krishnamoorthy, "An enhanced YOLOv8 for flame and smoke detection with dilated convolution and image dehazing," in Fourth International Conference on Signal Processing and Computer Science (SPCS 2023), 2023, vol. 12970: SPIE, pp. 604-608.

[19] K. Muhammad, J. Ahmad, and S. W. Baik, "Early fire detection using convolutional neural networks during surveillance for effective disaster management," Neurocomputing, vol. 288, pp. 30-42, 2018.

[20] Z. Jiao et al., "A deep learning based forest fire detection approach using UAV and YOLOv3," in 2019 1st International conference on industrial artificial intelligence (IAI), 2019: IEEE, pp. 1-5.

[21] S. T. Seydi, V. Saeidi, B. Kalantar, N. Ueda, and A. A. Halin, "Fire-Net: A deep learning framework for active forest fire detection," Journal of Sensors, vol. 2022, pp. 1-14, 2022.

[22] A. B. Abdusalomov, B. M. S. Islam, R. Nasimov, M. Mukhiddinov, and T. K. Whangbo, "An improved forest fire detection method based on the detectron2 model and a deep learning approach," Sensors, vol. 23, no. 3, p. 1512, 2023.

[23] C. Bahhar et al., "Wildfire and Smoke Detection Using Staged YOLO Model and Ensemble CNN," Electronics, vol. 12, no. 1, p. 228, 2023.

[24] F. M. Talaat and H. ZainEldin, "An improved fire detection approach based on YOLO-v8 for smart cities," Neural Computing and Applications, vol. 35, no. 28, pp. 20939-20954, 2023.

# Design and Analysis of Deep Learning Method for Fragmenting Brain Tissue in MRI Images

Ting Yang, Jiabao Sun*

School of Information and Electrical Engineering, Shaoxing University Yuanpei College, Shaoxing, Zhejiang, 312000, China

*Abstract*—An essential component of medical image processing is brain tumour segmentation. The process of giving each pixel a label is called image segmentation in order for pixels bearing the same label to share characteristics and help distinguish the target. A higher fatality rate and additional dangers can be avoided with early identification. It can be challenging and time-consuming to manually (man-made) segment brain tumours from the numerous MRI pictures generated during medical procedures in order to diagnose malignancy. This is the fundamental reason why brain tumour imaging has to be automated. The deep learning technique for the segmentation of brain tissue in magnetic resonance imaging (MRI) pictures was examined and enhanced in this work. Researchers are using deep learning techniques—convolutional neural networks in particular—to tackle the complex problem of biological image fragmentation object recognition. In contrast to traditional classification techniques that take in manually constructed qualities, convolutional neural networks automatically extract the required complicated features from the data itself. This solves a number of problems.

*Keywords—Brain tumor; deep learning; neural networks; magnetic resonance imaging*

## I. INTRODUCTION

Digital images are frequently utilised in image processing, and improved hospital healthcare services have been made possible in recent years by information technology and electronic health systems in the medical industry. When one or more malignant tissues in the brain grow abnormally, brain tumours result. There are two kinds of brain tumours: benign and malignant. A brain tumour is considered benign if the tumour tissues are growing uniformly and the cancer cells are dormant. A malignant tumour is one that spreads to all associated tissues if the cancerous tissues are non-uniform or if the cells are active, depending on the individual. Tumour location diagnosis is challenging due to the intricate structure and tissues of the human brain. Malignant tumours spread over the entire brain or spinal cord tissues, depending on when they are discovered. They do this by transferring diseased tissues to healthy tissues. Treatment for malignant tissues becomes more challenging and, in most circumstances, incurable, meaning the patient will eventually die from the disease as it spreads to additional regions. Thus, among the most important issues facing patients are early detection, type categorization, and infection rate.

With the development of new imaging techniques that make this information more accessible to doctors, radiation, surgery, or chemotherapy may be the best course of action. It follows that a patient's chances of surviving a tumour can be greatly raised if the tumour is accurately discovered in its early stages. Under the effect of imaging techniques, fragmentation is employed to identify the tumour area [1], [2]. Texture, contrast, border, and colour are the different components that make up a picture. As previously stated, aberrant and non-uniform growth of brain tissues is the definition of brain cancer or brain tumour [3], [4]. Brain tumours are among the deadliest forms of cancer, although being relatively uncommon. Brain tissue cells are the source of cells seen in primary brain tumours, and they begin in the brain. An MRI of the patient's brain is one of the best methods for diagnosing brain tumours. This technique facilitates a simple first diagnosis by giving medical professionals vital information on the structure, size, and metabolism of the brain tumour. This type of imaging is a common way to look at and diagnose brain tissue.

The study employed a variety of segmentation techniques, including a histogram, neural network segmentation methods, physical model-based techniques, clustering algorithms, Mean-Shift algorithm, k-means algorithm, Fuzzy C-Means Clustering (FCM) Algorithm, and Expect maximisation algorithms, to diagnose and classify the type and size of brain tumours in patients [5], [6]. Based on the segmentation system performance, various segmentation techniques are compared. For improved segmentation, artificial intelligence and enough information are typically combined [7], [8]. Conversely, the deep learning approach has shown the best results.

Abd-Ellah et al. reviewed the diagnosis of brain tumours using MRI scans in 2019 [9]. Additionally, the Support Vector Machine's (SVM) and Ecoc's error correction output codes are used to attain accuracy in the classification stage. PCA is used for feature extraction and selection in DWT machine learning. The collected features are classed using a seven-layer dynamic neural network (DNN) [10]. Terms and titles like CNN architecture, inverted graphics network, deep neural network, deep belief network, and so on are used in image processing technology; CNN architecture is the most commonly used word. Convolutional layers are typically used in the CNN architecture's feature extraction layer and input layer. The suggested method was proposed by Mohan et al. (2018) in three primary steps: CNN classification, post-processing, and pre-processing. Because CNN architecture performs so well in brain image recognition, medical professionals and researchers are using it more frequently than they used to because it is faster and more accurate than other methods. [11, 12]. A faultless network architecture is made possible by the CNN approach, which directly learns the difficult and complicated aspects of brain MRI images to aid in feature extraction and reduction. CNN receives brain MRI image excerpts as input,

and uses Neutral Density (ND) filters to extract complicated characteristics. The benefits of CNN architecture have been covered thus far; however, there are drawbacks to this approach as well, including high training computational costs, data consumption, and challenging hidden-layer complicated architecture design. It has also looked at the difficulties and issues that other studies have had when reporting the data, the type of tumour, and the algorithms' performance standards utilising the sensitivity, precision, and accuracy criteria. Based on the research and studies conducted, it is recommended that greater focus be placed on the identification, location, segmentation, and algorithm validation of brain tumours using magnetic resonance imaging.

When a brain tumour reaches an advanced stage, it can be difficult to diagnose it quickly. Magnetic resonance imaging is often chosen to improve the outcome of brain tumour diagnosis. It is exceedingly challenging to identify a tumour more accurately without harming healthy tissue. A new method for using magnetic resonance imaging (MRI) to diagnose brain tumours is presented in order to fix the flaws [13], [14]. As previously shown, early diagnosis of brain tumours can greatly aid patients in their recovery, contingent upon the nature and degree of metastasis. In the clinic or hospital, manually evaluating the numerous magnetic resonance imaging (MRI) scans that are generated on a regular basis is a challenging procedure. Thus, computer technologies should be employed for fast and accurate early diagnosis of brain tumours. Three key components are involved in diagnosing brain tumours using MRI images: segmentation, classification, and tumour diagnosis. Studies show that the majority of works in recent years have concentrated on the application of conventional learning machines. More precise techniques have been used recently by researchers; these are also used to find brain tumours. In general, this article's main theme gives a summary of both the development processes for new deep learning approaches for diagnosing brain cancer and traditional machine learning techniques.

The primary goals of the research presented in this paper are to provide an enhanced approach for brain segmentation, assess the effectiveness of the suggested approach, and review and compare the current approaches for brain tumour segmentation. The article's overall structure is set up so that the Section II provides a summary of the fundamental ideas and earlier approaches used in this area. The recommended dataset is introduced in the Section III, while the recommended approach is offered in the second portion. The steps for the study method are presented in the Section IV. The suggested approach is tested, its outcomes examined, and it is contrasted with alternative segmentation techniques in the Section V. The suggested algorithms are concluded and summarised in Section VI.

## II. PROPOSED METHODS

An enhanced strategy built on the deep learning technique is provided in this section. Compared to previous methods, the suggested method offers a higher sensitivity and accuracy. Neural networks can be used to implement the suggested method. Deep learning is one of the subgroups of machine learning. According to a set of techniques, deep learning

attempts to model high-level abstract notions in the data by using a deep graph with numerous processing levels made up of several linear and non-linear transformation layers [15], [16]. The study of neuronal behaviours in the human brain provides the learning structure with its primary purpose [17]. Models like deep neural networks, sophisticated neural networks, and deep belief networks have advanced well in the domains of image and natural language processing [18], [19]. In actuality, the study of novel artificial neural network techniques is referred to as deep learning. Deep learning can be used to meet the issues of avoiding the drawbacks of traditional machine learning methods [20], [21]. In several fields of medical image analysis, including computer-aided diagnosis, in-depth learning is becoming more and more common [22], [23]. As was previously noted, one of the subcategories of machine learning technology is the deep learning algorithm, which looks for many levels of distributed input data. The paper introduces a novel approach for segmenting tumour locations using a mix of artificial neural network and K-means fuzzy algorithm [37]. The process comprises of four stages, namely denoising, feature extraction and selection, classification, and segmentation. Initially, the preprocessed image is eliminated utilising the Wiener filter, subsequently followed by the extraction of significant GLCM features from the images. Subsequently, deep learning techniques were employed to categorize abnormal photos from normal images. Ultimately, the tumour region is segmented independently by subjecting it to the K-Means fuzzy algorithm. The proposed segmentation strategy is tested on the BRATS dataset and achieves an accuracy of 94%, a sensitivity of 98%, a specificity of 99%, and a Jaccard index of 96%.To address these issues, a fresh and enhanced deep learning model for image segmentation based on the cascaded regression (DLCR) technique is put forth. The fully convolutional neural network (FCNN) method is used in the suggested method to pre-process magnetic resonance imaging (MRI) pictures using the normalisation method. The data is then reduced and the matching feature from each feature vector is obtained through feature extraction using the Gaussian mixture model (GMM). Next, our suggested approach was contrasted with the existing techniques, which include the predictive machine learning model (MLPM), deep learning framework (DLF), and extreme learning machine local receiving fields (ELM-LRF). According to the findings, the suggested DLCR approach outperforms the current techniques in terms of sensitivity, specificity, recall ratio, precision ratio, peak signal-to-noise ratio (PSNR), and root mean square error (RMSE) [38].

These techniques are often divided into four groups: auto encoders, convolutional neural networks, sparse coding, and restricted Boltzmann machines (RBMS). Unlike human feature extraction in classical learning methods, machine learning technology may extract complex stages from their learning abilities. Through training, they get more acceptable findings in vast data sets. The fast rise in GPU processing capacity has enabled the creation of cutting-edge deep learning methods. Deep learning algorithms have been trained against picture alterations using millions of photos and resistance [24], [25].

An artificial neural network that resembles the structure of the human brain and is capable of carrying out mathematical

operations is called a deep learning algorithm. However, technological developments have given rise to methods for optimising traditional neural networks [26]. A multi-layer neural network with numerous hidden layers and free parameters does deep learning [27]. One of the most significant deep learning techniques is convolution neural networks, which effectively train several layers. This is one of the most popular and highly effective ways. The convolution layer, pulling layer, and fully connected layer are the three primary layers that make up a CNN network in general. Various layers carry out various functions. There are two steps in a convolutional neural network: the feed stage and the backpropagation stage for training.

Convolution is applied to each layer of the image after the point between the input and each neuron's parameters is multiplied in the first step. This step involves feeding the network input and then computing the network's output. The network error is computed and the network parameters are set using the output outcome. In order to calculate the error rate, compare the network output with the right response using a loss function. The post-publication stage starts in the following step, based on the computed mistake amount. Using the loss function, the network's output is compared to the right response to determine the error rate. The post-release stage starts in the following phase, depending on the quantity. This phase computes each parameter's gradient using the chain law. Every parameter has an effect on the network, according to the error made on each parameter based on the methodology. The feeding step is forward once the network has been modified and the parameters have been updated. Eventually, following numerous iterations of these procedures, the network's training is complete.

*A. CNN Network Layers*

Convolutional neural networks, in general, are neural networks with particular, hierarchical rules in which a number of fully connected layers come after the convolutional layers and the pulling layers. The convolution layer is one of the CNN layers. In these layers, the CNN network solves the middle feature map and the input picture using several cores, which results in distinct features in Fig. 1. There are three advantages to convolution operations: stability, immutability with regard to object relocation, local connectedness, and a significant reduction in the number of parameters.



Fig. 1.    Convolution layer performance.

The Pooling layer is the subsequent CNN layer. Pooling layers, which are frequently placed after a convolutional layer,

can be used to reduce the size of the feature map and network parameters. Convolution layers and other pooling layers are resilient against displacement because they incorporate nearby pixels into their computations. The most popular pooling layer implementations make use of the max (max pooling) and average (average pooling) functions. It accelerates convergence, enhances generalisation, and produces a well-chosen set of fixed features. Fig. 2 displays a max-pooling result.



Fig. 2.    Max-pooling performance.

The fully linked layer is the final layer of CNN. As seen in Fig. 3, there are fully linked layers following the final polishing layer, each of which is connected to the neurons in the layer before it. On the other hand, about 90% of the parameters of a CNN network are made up of fully linked layers, which are likewise conventional.



Fig. 3.    Fully connected layers performance.

This kind of layer's main problem is its excessive number of parameters. As a result, there is a significant processing cost that needs to be covered by training. Therefore, cutting back on connections and eliminating these layers entirely using various techniques is a popular technique that yields good results.

*B. Enhanced Neural Network Capabilities*

The advantage of deep learning over surface learning is the ability to build deep architectures with higher data volumes and higher abstract information transmission. However, overfitting may become a problem if a lot of new parameters are included. A plethora of regularisation techniques have been introduced recently to address overfitting and enhance network

performance. Note that not every method covered in this section is unique to neural networks. This implies that various methods, such as the dropout method, can be used with conventional artificial neural networks.

*1) Deep learning dropout:* The purpose of this approach is to prevent over-coverage. Each neuron in the network is either retained with a probability of p or eliminated with a probability of p-1 at each training stage. Eventually, just a smaller network is left. An expelled node's input and output edges are also eliminated. At that point, the data will only be used to train the smaller network. Next, the removed nodes are removed from the network and then added back in with their original weights. A comparison between the networks is shown in Fig. 4.

Drop Connect is a well-liked variation of Dropout that drops weights at random rather than activation values. Research has demonstrated that, albeit more slowly, this approach can regularly outperform the Dropout method in a variety of common benchmark categories.

*2) Deep clustering in deep learning:* Without any prior knowledge of different factors like feature selection, distance criteria, or clustering techniques, the goal of clustering is to group comparable data points. One deep learning algorithm is deep clustering. Two-dimensional non-linear data representations from the data set are employed for learning. Quantization of deep neural networks requires network loss. Deep clustering techniques now in use can be generally classified into two groups. After learning a representation that jointly optimises learning and clustering, a two-step challenge is implemented. The first class of algorithms makes use of sophisticated unsupervised learning frameworks and methodologies directly. An additional set of algorithms simulates a classification error in monitoring by attempting to precisely characterise a clustering failure.

*3) Data redundancy in deep learning:* The utilisation of redundant data results from the redundancy of data generated when CNN is utilised for object detection. Forecasts may benefit from the addition of new data, which also improves the data's quality and accuracy. This can be used to decrease interference and expand training set sizes. [28, 29]

*4) In deep learning, auto encoder and deep auto encoder:* A neural network that has been trained to replicate its input to its output is called an autoencoder. Neural networks whose input and output are identical are an artificial neural network type that is used to encode effective learning data in an unsupervised manner. They function by first compressing the input into a representation of hidden space, from which the output is then rebuilt. An autoencoder aims to encrypt a dataset, usually with the purpose of reducing dimensionality by conditioning the network to reject spurious signals. Fig. 5 depicts an autoencoder's full workflow. You train an autoencoder to re-encode your input X, as opposed to training the network to anticipate the target value Y for the input X. During this procedure, the auto-encoder is optimised by minimising the update error.

Hinton [28] introduced the deep autoencoder, which is still being researched extensively in current studies. As previously noted, a backpropagation process, like the gradient approach, is used to train the autoencoder. Ignoring the benefits, this model's significant drawback becomes apparent when an error happens, and these layer flaws can cause the model to become extremely inefficient. The network reconstructs the training data average as a result of this error. This problem can be effectively solved by pre-training the network with initial weights, which comes close to the ultimate solution. Simultaneously, several auto-encoders are available that promise to retain representations for an extended period of time due to continuous changes in input. The autoencoder is compelled by full representation learning to extract the most important features from the data.



Fig. 4. Comparison between the networks.

Fig. 5. The autoencoder's general process.



Fig. 6. Denoising autoencoder.

*5) Denoising autoencoder in deep learning:* In Fig. 6, the DAE procedure is displayed. (DAE) is a power-boosting model known as an automatic denoising encoder. Vincent presents this model, which is able to discern between the accurate input and the tampered version. Put another way, have the model see the input distribution's structure.

In actuality, this task's objective is to either eliminate or clean up the contaminated input. The depiction of higher levels that are comparatively resistant to corrupt inputs and the extraction of characteristics with a relevant structure in the input distribution are the two linked hypotheses of this technique. Its ability to correctly recover from a damaged version is one of its excellent features, and it is noise-resistant.

*C. Network in Network (NIN) and Activation Functions*

Fully linked layers are used for image processing after a series of layers and aggregation layers are used to extract features from spatial structures. The primary idea behind it is to substitute a perceptron neural network (many fully connected layers with non-linear activation functions) for the

convolutional layer and learn its parameters from the data. In this sense, non-linear neural networks take the place of linear filters. This technique produces superior picture categorization results. It is in charge of the data's non-linear transformation. The modified linear unit (ReLU) is computed in Eq. (1).

$$F(x)=max\ (0,\ x) \qquad (1)$$

It has been discovered that these functions outperform hyperbolic tangent functions or classical functions in terms of training speed [30], [31]. On the other hand, applying a constant 0 can harm the flow gradient and the subsequent weight adjustment [32], [33]. Using a variable named Leaky Rectified Linear Unit (LReLU), which adds a tiny slope to the performance's negative side, we adjust to these constraints. Eq. (2) is used to calculate this function.

$$F(x)=max\ (\alpha x,\ x) \qquad (2)$$

The activation function is among the most often used (ReLU). The ReLU non-linear unit's job is to clip any negative input value in the data to zero, reducing it, and then taking the positive input values as the output [34], [35].

## III. DATA SET

We utilised the publically accessible brain tumour dataset from Figshare [4] to examine and assess our suggested methodology, employing various convolutional neural network designs. Cheng developed it in 2017. The dataset consists of 3064 brain MRI slices obtained from 233 people. It encompasses three types of brain tumours:

The total amount of images for meningioma is 708, for pituitary it is 930, and for glioma tumour it is 1426. The dataset is accessible to the public via the Figshare website in the ".mat" format, which is compatible with MATLAB. Each MAT-file consists of a structure that includes a patient ID, a distinctive label indicating the type of brain tumour, picture data in a 512 × 512 format represented as unsigned 16-bit integers, a vector providing the coordinates of discrete points that define the tumour perimeter, and a binary mask image representing the ground truth. Fig. 7 shows the eight examples of database images.

The enhanced deep learning algorithm's pseudo-code is presented in the text below.

An overview of CNN architecture is presented in Fig. 8.

Algorithm 1: Improved Deep Embedded Clustering

***Input:** Input data: X; Number of clusters: K; Target Distribution update interval: T; Stopping Threshold: $\delta$ ; Maximum iterations: MaxIter.*
*Output: Autoencoder's weights W and W' ; Cluster centers $\mu$ and labels s.*

**1** *Initialize $\mu$, W' and W according to Section 3.1.*
**2** *for iter $\in \{0,1,\ldots,MaxIter\}$ do*
**3** *if iter %T == 0 then*
**4** *Compute all embedded points $\{z_{i=} f_w(xi)\}i = 1$*
**5** *update P using (3), (4) and {zi} i=1.*
**6** *Save last label assignment: $S_{old}$=s.*
**7** *Compute new label assignment s nia (14).*
**8.** *if sum $(S_{old} \neq s)/n < \delta$ then*
**9.** *Stop training.*
**10** *Choose a batch of sample S $\in$ X.*
**11** *Update $\mu$, W' and W via (11), (12) and (13) on S.*

## III. DATA SET

We utilised the publically accessible brain tumour dataset from Figshare [4] to examine and assess our suggested methodology, employing various convolutional neural network designs. Cheng developed it in 2017. The dataset consists of 3064 brain MRI slices obtained from 233 people. It encompasses three types of brain tumours:

The total amount of images for meningioma is 708, for pituitary it is 930, and for glioma tumour it is 1426. The dataset is accessible to the public via the Figshare website in the ".mat" format, which is compatible with MATLAB. Each MAT-file consists of a structure that includes a patient ID, a distinctive label indicating the type of brain tumour, picture data in a 512 × 512 format represented as unsigned 16-bit integers, a vector providing the coordinates of discrete points that define the tumour perimeter, and a binary mask image representing the ground truth. Fig. 7 shows the eight examples of database images.

The enhanced deep learning algorithm's pseudo-code is presented in the text below.

An overview of CNN architecture is presented in Fig. 8.

Algorithm 1: Improved Deep Embedded Clustering

***Input:** Input data: X; Number of clusters: K; Target Distribution update interval: T; Stopping Threshold: $\delta$ ; Maximum iterations: MaxIter.*
*Output: Autoencoder's weights W and W' ; Cluster centers $\mu$ and labels s.*

**1** *Initialize $\mu$, W' and W according to Section 3.1.*
**2** *for iter $\in \{0,1,\ldots,MaxIter\}$ do*
**3** *if iter %T == 0 then*
**4** *Compute all embedded points $\{z_{i=} f_w(xi)\}i = 1$*
**5** *update P using (3), (4) and {zi} i=1.*
**6** *Save last label assignment: $S_{old}$=s.*
**7** *Compute new label assignment s nia (14).*
**8.** *if sum $(S_{old} \neq s)/n < \delta$ then*
**9.** *Stop training.*
**10** *Choose a batch of sample S $\in$ X.*
**11** *Update $\mu$, W' and W via (11), (12) and (13) on S.*

Fig. 7.    The eight examples of database images.



Fig. 8.    CNN architecture.

## IV.    STUDY METHOD STEPS

Early diagnosis, therapy, and follow-up are made possible by having an effective brain tumour detected by multimodal imaging as soon as possible [36]. In some social networks, deep learning algorithms have quickly taken the lead in medical picture analysis. We investigate deep learning applications in object detection, segmentation, recording, and image classification. The proposed deep learning models work well for brain tumour segmentation and provide very accurate results. The suggested models can also assist medical practitioners in shortening the time needed for diagnosis. The goal of this research is to create a fully automated model for MRI brain tumour segmentation that operates without the need for human intervention. Instead of having multiple boxes (or networks) from input to output, a single grid is used in this deep learning model's proposal. Deep-end learning typically uses a neural network in place of the numerous processing stages that standard machine learning techniques typically call for. In order to get the most out of each layer before going on to the next, we construct this final model. When treating malignant tumours, deep learning models are employed. These tumours can arise anywhere in the brain and have four unpredictable characteristics: varying sizes, forms, and contrast. Therefore, convolutional neural networks are the best option for machine learning.

### A.  Pre-processing

To get better results, adjustments to the photos are required. Understanding, processing, and picture analysis are some of these adjustments. These modifications enhance image processing systems to enable them to carry out tasks more quickly and accurately. Pre-processing, picture quality enhancement, image conversion, categorization, and image analysis are the primary procedures that raise the efficiency of these systems, in that order. These techniques involve analysing images for certain goals, with computer-generated mathematical equations simulating irregular human visual features. Computer vision is the scientific image analysis used in many areas of engineering, health, imaging, security, and astronautics. Modern digital technology is able to operate multidimensional switches and many parallel computers thanks to systems derived from simple digital circuits.

One of the things like image processing, picture analysis, and image perception is the goal of these modifications. Since the majority of remote sensors store their data digitally, digital processing is eventually needed for all picture interpretation and analysis. In order to allow further interpretation and analysis, digital image processing may comprise a variety of processes, such as alteration, digital optimisation, data formatting, or even computer goals automation and automatic properties. To carry out this procedure, the data must be in a format that is appropriate for physical storage.

Four categories comprise the most often utilised processing procedures in image analysis systems: pre-processing, image enhancement, image transformation, image analysis, and classification. These procedures fall under the general category of radiometric or geometric adjustments and are typically grouped prior to the primary picture analysis and information extraction. To get an accurate picture of the radiation that the receivers are receiving, radiometric corrections entail adjusting data for undesired noise and anomalies in the receivers. The purpose of geometric corrections is to simplify the image and translate it into actual coordinates on the surface of the earth.

The procedures used in the Image Enhancement category are only intended to raise and enhance the image's resolution in order to provide a better understanding of the image. These acts are hypothetically comparable to the preceding category in picture transmission. However, this group involves the combined processing of data acquired from various spectral bands, in contrast to this group, which pertains to only one data channel. The original bands are combined with mathematical operations (addition, subtraction, multiplication, and division) to create new images that are more readable or have better qualities. Pixel classification and analysis aim to identify and characterize data pixels. The process of classifying assigns each pixel in an image to a group or theme based on the statistical characteristics of its floating values. It is typically applied to multi-channel data groups. There are two main approaches to this: supervisor-free and supervisor-involved.

### B. Image Transformation Color to Gray

Three colors green, blue, and red combine to form the most common colour model in computer graphics. These three elements can combine to give a total of 16581375 distinct colours. In computer graphics, this colour model is referred to as RGB. Apart from the RGB colour model, there exist various additional models that display colours differently, including CMYK, HSI, HSV, and Grayscale. We are particularly interested in the Grayscale colour model in the interim. Because a grayscale image will do in the majority of applications, and a colour image is not necessary. Generally speaking, the grey image is referred to as the black and white image (of course, calling something grey instead of black and white is incorrect; this is just being clarified). Three matrices, one for each of the colours red, green, and blue in the image, make up an RGB image. Three matrices are created by mixing the values from the appropriate verses to display the image on the screen. For the majority of applications, a grayscale image will do; a colour image is not necessary. When the values of the R, G, and B components of a pixel are the same, the pixel will have a grey value. In accordance with this concept, he utilised Eq. (3) to grayscale input images that are RGB.

$$S\_R(x,y) = S\_G(x,y) = S\_B(x,y) = \frac{[R(x,y) + G(x,y) + B(x,y)]}{3} \quad (3)$$

where, R, G, and B are the matrices representing the red, green, and blue components of the input image, and S_X are the components of the output image.

### C. Median Filter

Among the non-linear filters is the median filter. Signals and noise in images are captured using this filter. Picture noise refers to background and other picture detections. For instance, you have to take a picture using one of the filters like the median noise filter—before you can identify a corner. When processing images, the median filter is frequently used. It is also impossible to overlook the median filter's application in signal processing. The median filter is utilised in blood pressure monitors, EEG systems, and radiography systems, among other specialised applications. We can list median and fashion filters as examples of non-linear filter types. The median filter is a low-pass filter that uses a face neighbourhood to sort the neighbourhoods in ascending order before sorting the middle element of the numbers to replace the centre pixel. It should be mentioned that noise from salt and pepper can be eliminated using the mid-pass filter.

### D. SOBLE Filter

There are three ways to convert a colour image to grayscale: edge features, edge detection, and image edging techniques. As using image edging techniques, the light intensity dramatically changes as points are marked in the image for edge identification. Changes in image properties are indicative of major changes in environmental factors. Image processing and feature extraction researchers are studying edge identification. The boundary in edge detection is found between regions of the grey surface with comparatively distinct characteristics. The fundamental idea behind the majority of edge detection techniques is the computation of a local derivative operator. It should be noticed at this point that the edge—the change from dark to light—is modelled as a gradual, rather than abrupt, grey surface change. This model demonstrates how sampling typically causes the borders of digital images to become slightly blurry.

### E. Image Binary Threshold

A threshold value must be used to establish the threshold when converting a grayscale image to a binary image. Pixels that are less than the threshold are assigned a value of 0, while those that are less than the threshold are assigned a value of 1 or 255. Thresholding an image can be done in various ways. They can be broadly separated into the five sections listed below. 1) Histogram-based techniques have been examined, including the analysis of the peaks, valleys, and curvatures of the smoothed histogram. 2) Clustering-based techniques, in which a grayscale image is divided into foreground and background. 3) Entropy-based techniques: these techniques determine the threshold based on the intersection of the two classes in the image as well as the background and background entropy distributions. 4) Object property-based methods: these determine the appropriate threshold by calculating the degree of similarity between the image and a binary and grayscale scale. 5) Location-based techniques: the desired threshold is computed using a higher-order probability correlation or distribution between pixels.

### F. Tumor Detection

Automatically detecting and extracting a portion of a tumor from a brain image is a challenging and complex mission for doctors and physicians.

## V. DISCUSSION

Systems for detecting brain tumours are employed in numerous sectors. One of the main objectives of this field's study is to use deep learning to develop a solution that complies with the guidelines and standards of brain tumour diagnosis systems, hence improving the precision of medical specialists. MATLAB software is used to perform the simulations. Its accuracy, precision, and sensitivity are higher than those of the approaches offered, based on the evaluation findings. The research methodology was covered in the previous section; the simulation findings are reported in this section. Fig. 9 shows the original image (a) and the image with a balanced histogram (b). Histogram balancing is the initial stage of image processing that enhances image quality.

The histogram in Fig. 10(a) shows the image prior to balancing and in Fig. 10(b) is the image following balancing. The horizontal axis of the histogram represents the colours, while the vertical axis represents the quantity of pixels in each colour. When using image histogram balancing, the image extends between 0 and 255 if the colour is between 0 and 100 when balanced.



Fig. 9.  (a) Original image, (b) Histogram balancing.



Fig. 10.  Image histogram, (a) Original, (b) Equalized.

In Fig. 11, the median filter operation is performed on the main image. To improve those pixels that is not of good quality. It is actually to improve the quality of images.

Fig. 12 shows the Sobel image and the original image is binaries in Fig. 13. In order to isolate the objects in a digital image from their backdrop and analyse them, the image is first converted to a binary image. This produces a binary image that is less in volume than the original image.

In the part of the images that are interconnected, the box is drawn with red lines around it by writing the program, as shown in Fig. 14. In other words, it actually features extraction.

Fig. 11. Median filter.



Fig. 12. Sobel image.



Fig. 13. Binarization.



Fig. 14. Connected component.

### A. Tumour Detection

A comparison of the results with other approaches is made, taking into account various circumstances to demonstrate the effectiveness of the suggested method. Simulated approaches are identical in terms of conditions and evaluation parameters. Lastly, graphs are used to plot and analyse the simulation findings.

Each of the 766 data points is put in for loop in the MATLAB programme with the intention of being worked on. Following the use of the enhanced DEEP LEARNING algorithm to identify tumour sites in the output photos, these 14 sample images are displayed. Brain tumours are depicted in these pictures as pink in Fig. 15.

### B. Evaluation Results

There are four sections to the evaluation results. 1) Positives that are accurately identified as true positives (TP). 2) False positives (FP), or negatives with a valid diagnosis. 3) Mistaken positives that are actually true negatives (TN). 4) Misdiagnosed negatives, or false negatives (FN). The assessed positive characteristics for simulating the first through sixth scenarios using various techniques are displayed in Table I.

Table II shows the evaluated negative parameters to simulate the seventh to twelfth scenarios for different methods.

Fig. 15. Tumor detection results.

TABLE I.    POSITIVE EVALUATION PARAMETERS

| Scenario | Method | Accuracy (%) | TP Ratio (%) | Precision TP (%) | Recall TP (%) | F-Measure (%) | Precision (%) |
|---|---|---|---|---|---|---|---|
| 1 | proposed method<br>Neural network<br>SVM<br>correlation<br>Euclidean | 95<br>90<br>79<br>74<br>58 | - | - | - | - | - |
| 2 | proposed method<br>Neural network<br>SVM<br>correlation<br>Euclidean<br>differential<br>block (8×8) differential | - | 98<br>96<br>95<br>94<br>87<br>48<br>40 | - | - | - | - |
| 3 | proposed method<br>Neural network<br>SVM<br>correlation<br>Euclidean | - | - | 93<br>74<br>69<br>66<br>63 | - | - | - |
| 4 | proposed method<br>Neural network<br>SVM<br>correlation<br>Euclidean<br>differential | - | - | - | 98<br>96<br>95<br>94<br>48<br>40 | - | - |
| 5 | proposed method<br>Neural network<br>SVM<br>correlation | - | - | - | - | 95<br>90<br>79<br>74 | - |
| 6 | proposed method<br>Neural network<br>SVM | - | - | - | - | - | 95<br>90<br>40 |

TABLE II.    NEGATIVE EVALUATION PARAMETERS

| Scenario | Method | Accuracy Total (%) | TN Ratio (%) | Precision TN (%) | Recall TN (%) | Accuracy TN (%) | F-Measure TN (%) |
|---|---|---|---|---|---|---|---|
| 1 | proposed method<br>Neural network<br>SVM<br>correlation | 93<br>69<br>66<br>63 | - | - | - | - | - |
| 2 | proposed method<br>Neural network<br>SVM<br>correlation | - | 60<br>52<br>37<br>30 | - | - | - | - |
| 3 | proposed method<br>Neural network<br>SVM<br>correlation<br>Euclidean<br>differential<br>block (8×8) differential | - | - | 66<br>62<br>60<br>58<br>55<br>24<br>22 | - | - | - |
| 4 | proposed method<br>Neural network<br>SVM<br>correlation | - | - | - | 60<br>52<br>49<br>46 | - | - |
| 5 | proposed method<br>Neural network<br>SVM<br>correlation<br>Euclidean<br>differential | - | - | - | - | 66<br>62<br>60<br>58<br>24<br>22 | - |
| 6 | proposed method<br>Neural network<br>SVM | - | - | - | - | - | 60<br>52<br>24 |

TABLE III.    RELATIONSHIPS RELATED TO THE EVALUATION PARAMETERS

| - | Reference Standard Stroke | Reference Standard No Stroke | - |
|---|---|---|---|
| **Self-report Stroke** | (TP) | (FP) | $PPV=TP/TP+FP$ |
| **Self-report No Stroke** | (FN) | (TN) | $NPV=TN/FN+TN$ |
| - | Sensitivity=$TP/TP+FN$ | Specificity=$TN/FP+TN$ | - |

In Table III, the relationships related to the evaluation parameters are calculated.

*1) Positive evaluation results:* The accuracy diagram, which compares the suggested method with the neural network, SVM, correlation, and Euclidean methods in Fig. 16, shows the green colour of the suggested method, the red colour of Euclidean, the black colour of correlation, the blue colour of the support vector machine, and the purple colour of the neural network. The training data percentage in the data set is shown by the horizontal axis in this image, and the accuracy %, which ranges from 0 to 1, is represented by the vertical axis. The diagram indicates that for the majority of the data, the suggested approach is more accurate than the other methods in every evaluation point.

The proposed method has been compared with SVM, Neural Network, Correlation, Euclidean, Differential, and Block (8x8) Differential methods in Fig. 17. The second scenario is the TP ratio diagram, which is the green colour of the proposed method, the red colour of Euclidean, the black colour of correlation, the blue colour of the support vector machine, the purple colour of the neural network, the blue colour of block 8x8, and the black colour of differential. In this case, the vertical axis is a percentage, while the horizontal axis represents the trained data % in the dataset. The suggested technique outperforms the other methods in all evaluation points for the majority of the data, as shown by the diagram.

The proposed method is compared with SVM, Neural Network, Correlation, and Euclidean methods in Fig. 18. The third scenario is the precision diagram, where the green colour represents the suggested method, the red colour represents the Euclidean method, the black colour represents a correlation, the blue colour represents a support vector machine, and the purple colour represents a neural network. In this case, precision is the vertical axis and the trained data percentage in the dataset is the horizontal axis. The diagram indicates that for the majority of the data, the recommended approach has outperformed alternative methods in every evaluation point.

The F-measure diagram, which compares the suggested method with the SVM, Neural Network, and Correlation methods in Fig. 20, represents the fifth scenario and Fig. 19 shows the recall diagram. It is coloured red for the proposed method, red for the correlation, blue for the support vector machine, and black for the neural network. The trained data % in the dataset is the horizontal axis in this case, and the f-measure percentage rate is the vertical axis. The diagram indicates that for the majority of the data, the recommended approach has a higher f-measure than the other methods in every evaluation point.


Fig. 16. Accuracy diagram.


Fig. 17. TP Ratio.


Fig. 18. Precision diagram.

Fig. 19. Recall diagram.


Fig. 20. F-measure diagram.

Correlation methods, as shown in Fig. 23. The horizontal axis represents the percentage of trained data in the data set, while the vertical axis represents the ratio of true negatives (TN). Based on the diagram, the suggested technique exhibits a higher TN Ratio compared to other methods at all evaluation points for the majority of the data.


Fig. 21. Precision total diagram.


Fig. 22. Accuracy total diagram.

The sixth scenario, which contrasts the suggested technique with SVM and neural network methods in Fig. 21, shows the precision total diagram with the red colour of the suggested method, the black colour of the neural network, and the blue colour of the support vector machine.

In this case, the vertical axis represents the overall accuracy % while the horizontal axis represents the percentage of trained data in the dataset. The diagram indicates that for the majority of the data, the recommended approach has outperformed alternative methods in every evaluation point.

*2) Negatives evaluation results:* The seven sonorities are the total accuracy chart, which contrasts the red, green, blue, and black colours of the neural network, support vector machine, and SVM in the proposed method with the neural network and correlation methods in Fig. 22. The training data % in the dataset is shown by the horizontal axis in this example, while the total correctness percentage is represented by the vertical axis. The diagram indicates that for the majority of the data, the suggested approach is more accurate than alternative methods in all evaluation points.

Scenario eight involves the TN Ratio chart, which displays a comparison between the green colour representing the proposed method, the red colour representing correlation, the blue colour representing the neural network, and the black colour representing the support vector machine. This comparison is made with the SVM, Neural Network, and


Fig. 23. TN Ratio diagram.

The precision TN diagram, which is coloured as follows: green for the suggested approach, red for Euclidean, black for the differential, blue for the support vector machine, and purple for the neural network, represents the ninth case. In Fig. 24, the suggested approach correlation is contrasted with the Differential, SVM, Euclidean, Neural Network, Correlation, and Block (8×8) Differential methods, as indicated by the black colour. In this case, the vertical axis represents the

percentage of TN accuracy, and the horizontal axis represents the proportion of trained data in the dataset. The suggested method outperforms competing methods in all evaluation points for the majority of the data, as shown by the chart.



Fig. 24. Precision TN diagram.



Fig. 25. Recall the TN diagram.

In Fig. 25, the Recall TN diagram presents a comparison between the proposed method and other methods, namely Euclidean, correlation, support vector machine (SVM), and neural network. The proposed method is represented by the green colour, Euclidean by red, correlation by black, SVM by blue, and neural network by purple. The horizontal axis represents the percentage of training data in the dataset, while the vertical axis represents the percentage of Recall TN. Based on the chart, the suggested technique exhibits higher Recall TN values compared to other methods across all evaluation points for the majority of the data.

The 11th scenario presents the TN accuracy diagram, where the proposed method is represented by black, Euclidean by red, correlation by black, support vector machine by red, neural network by blue, and differential by yellow. This diagram compares the proposed method with SVM, Neural Network, Correlation, and Euclidean methods. Fig. 26 presents a comparison of the differential. The horizontal axis represents the percentage of learned data in the dataset, while the vertical axis represents the accuracy rate of True Negative (TN). Based on the diagram, the suggested method demonstrates superior accuracy in true negatives (TN) compared to other methods across all evaluation points for the majority of the data.

Scenario twelve involves the comparison of the proposed technique, support vector machine (SVM), and neural network

using the F-measure TN diagram. The proposed method is represented by the colour purple, SVM by black and neural network by red. This comparison is depicted in Fig. 27. The horizontal axis in this scenario represents the percentage of trained data in the dataset, while the vertical axis represents the F-Measure TN. Based on the diagram; the suggested method consistently outperforms other methods in terms of F-Measure TN across all evaluation points for the majority of the data.

The performance of the suggested technique has been evaluated by comparing it with SVM, Neural Network, Correlation, Euclidean, Differential, and Block (8×8) Differential methods. The techniques are executed within the MATLAB simulation environment. The criteria of accuracy, precision, and sensitivity have been chosen for comparison and suggested for evaluating the approaches. The evaluation findings have been graphically shown and demonstrate that the suggested method outperforms previous methods in terms of accuracy, precision, sensitivity, and recall in all scenarios.

Table IV presents a comparison of distinct datasets that utilise various approaches with fused segmented images. The image datasets undergo pre-processing using image fusion methodology and are subsequently segmented using U-Net. Conversely, the method of merging images is not employed; just the segmentation of images is performed. It is evident that the techniques employed with fused images yield output that is comparable to that of non-fused images. The image fusion approach involves combining images obtained from diverse sensors, which contain essential data, using various numerical models to get a unified composite image.



Fig. 26. Accuracy TN diagram.



Fig. 27. F-measure TN diagram.

TABLE IV. COMPARATIVE ANALYSIS OF DIFFERENT PARAMETERS IN TERMS OF BRAIN TISSUE REGION IDENTIFICATION AND CLASSIFICATION USING MULTIMODAL IMAGING METHODS

| Image Fusion | Image Segmentation Method | Classification Method | Dataset | Data Volume | Data Class | Accuracy (%) |
|---|---|---|---|---|---|---|
| Non-Fused Image | U-Net | U-Net | Brain tumour fig share | 3074 | 3 | 96.21 |
| Fused Image | U-Net | ResNet50 | BraTS 2012 | 3074 | 3 | 95.34 |
| Fused Image | U-Net | U-Net | BraTS 2018 | 3074 | 3 | 89.37 |
| Fused Image | U-Net | CNN | BraTS 2018 | 3074 | 3 | 90.01 |
| Non-Fused Image | U-Net | U-Net | BraTS 2012 | 3074 | 3 | 97.68 |
| Fused Image | U-Net | ResNet50 | BraTS 2012 | 3074 | 3 | 88.69 |
| | | | Dataset Figshare | 3074 | 3 | 98.71 |

## VI. CONCLUSION

This article presents a neural network-based method for automatically detecting brain tumours. Various methodologies are being examined. The evaluation method presents the results based on accuracy, precision, and sensitivity. The initial step involves presenting the pre-processing outcomes, encompassing histogram equalisation, median filtering, edge detection, binarization, and identification of related pixels. Subsequently, the detection outcomes for 766 data points are presented, with 14 of them being displayed. Subsequently, the evaluation relationships are articulated with respect to accuracy, correctness, and evaluation standards. Ultimately, the evaluation outcomes for the positives and the evaluation outcomes for the negatives are disclosed. The simulation is conducted with 766 data points, and the assessment criteria are also assessed. The proposed strategy for improvement is the utilisation of deep learning, which has enhanced the effectiveness of the previous method. The accuracy of these is verified, and subsequently, the detection approach is established and compared to alternative methodologies. In future research, it is possible to employ optimisation methods like as genetic algorithms and colonial competition, among others, in the field of deep learning as alternatives to the modified neural network.

## REFERENCES

[1] T. Logeswari and M. Karnan, "An improved implementation of brain tumor detection using segmentation based on hierarchical self-organizing map," International Journal of Computer Theory and Engineering, vol. 2, no. 4, p. 591, 2010.

[2] S. Agrawal, R. Panda, and L. Dora, "A study on fuzzy clustering for magnetic resonance brain image segmentation using soft computing approaches," Appl Soft Comput, vol. 24, pp. 522–533, 2014.

[3] N. Gordillo, E. Montseny, and P. Sobrevilla, "State of the art survey on MRI brain tumor segmentation," Magn Reson Imaging, vol. 31, no. 8, pp. 1426–1438, 2013.

[4] J. Liu, M. Li, J. Wang, F. Wu, T. Liu, and Y. Pan, "A survey of MRI-based brain tumor segmentation methods," Tsinghua Sci Technol, vol. 19, no. 6, pp. 578–595, 2014.

[5] M. Prastawa, E. Bullitt, S. Ho, and G. Gerig, "A brain tumor segmentation framework based on outlier detection," Med Image Anal, vol. 8, no. 3, pp. 275–283, 2004.

[6] N. Moon, E. Bullitt, K. Van Leemput, and G. Gerig, "Automatic brain and tumor segmentation," in Medical Image Computing and Computer-Assisted Intervention—MICCAI 2002: 5th International Conference Tokyo, Japan, September 25–28, 2002 Proceedings, Part I 5, Springer, 2002, pp. 372–379.

[7] B. H. Menze, K. Van Leemput, D. Lashkari, M.-A. Weber, N. Ayache, and P. Golland, "A generative model for brain tumor segmentation in multi-modal images," in Medical Image Computing and Computer-Assisted Intervention–MICCAI 2010: 13th International Conference, Beijing, China, September 20-24, 2010, Proceedings, Part II 13, Springer, 2010, pp. 151–159.

[8] D. Kwon, R. T. Shinohara, H. Akbari, and C. Davatzikos, "Combining generative models for multifocal glioma segmentation and registration," in Medical Image Computing and Computer-Assisted Intervention–MICCAI 2014: 17th International Conference, Boston, MA, USA, September 14-18, 2014, Proceedings, Part I 17, Springer, 2014, pp. 763–770.

[9] M. K. Abd-Ellah, A. I. Awad, A. A. M. Khalaf, and H. F. A. Hamed, "A review on brain tumor diagnosis from MRI images: Practical implications, key achievements, and lessons learned," Magn Reson Imaging, vol. 61, pp. 300–318, 2019.

[10] Zhao, Y., Liang, H., Zong, G., & Wang, H. (2023). Event-Based Distributed Finite-Horizon $ H\_\infty $ Consensus Control for Constrained Nonlinear Multiagent Systems. IEEE Systems Journal.

[11] G. Mohan and M. M. Subashini, "MRI based medical image analysis: Survey on brain tumor grade classification," Biomed Signal Process Control, vol. 39, pp. 139–161, 2018.

[12] Blbas, H., & Kadir, D. H. (2019). An application of factor analysis to identify the most effective reasons that university students hate to read books. International Journal of Innovation, Creativity and Change, 6(2), 251-265.

[13] Samiei, M., Hassani, A., Sarspy, S., Komari, I. E., Trik, M., & Hassanpour, F. (2023). Classification of skin cancer stages using a AHP fuzzy technique within the context of big data healthcare. Journal of Cancer Research and Clinical Oncology, 1-15.

[14] Cheng, F., Niu, B., Xu, N., Zhao, X., & Ahmad, A. M. (2023). Fault detection and performance recovery design with deferred actuator replacement via a low-computation method. IEEE Transactions on Automation Science and Engineering.

[15] Cao, C., Wang, J., Kwok, D., Cui, F., Zhang, Z., Zhao, D., ... & Zou, Q. (2022). webTWAS: a resource for disease candidate susceptibility genes identified by transcriptome-wide association study. Nucleic acids research, 50(D1), D1123-D1130.

[16] N. Le Roux and Y. Bengio, "Representational power of restricted Boltzmann machines and deep belief networks," Neural Comput, vol. 20, no. 6, pp. 1631–1649, 2008.

[17] Xiao, L., Cao, Y., Gai, Y., Khezri, E., Liu, J., & Yang, M. (2023). Recognizing sports activities from video frames using deformable convolution and adaptive multiscale features. Journal of Cloud Computing, 12(1), 1-20.

[18] Sajadi, S. M., Kadir, D. H., Balaky, S. M., & Perot, E. M. (2021). An Eco-friendly nanocatalyst for removal of some poisonous environmental pollutions and statistically evaluation of its performance. Surfaces and Interfaces, 23, 100908.

[19] A. Hyvärinen and P. Hoyer, "Emergence of phase-and shift-invariant features by decomposition of natural images into independent feature subspaces," Neural Comput, vol. 12, no. 7, pp. 1705–1720, 2000.

[20] Sun, J., Zhang, Y., & Trik, M. (2022). PBPHS: a profile-based predictive handover strategy for 5G networks. Cybernetics and Systems, 1-22.

[21] Trik, M., Akhavan, H., Bidgoli, A. M., Molk, A. M. N. G., Vashani, H., & Mozaffari, S. P. (2023). A new adaptive selection strategy for reducing latency in networks on chip. Integration, 89, 9-24.

[22] Omer, A. W., Blbas, H. T., & Kadir, D. H. (2021). A Comparison between Brown's and Holt's Double Exponential Smoothing for Forecasting Applied Generation Electrical Energies in Kurdistan Region.

[23] Wang, G., Wu, J., & Trik, M. (2023). A Novel Approach to Reduce Video Traffic Based on Understanding User Demand and D2D Communication in 5G Networks. IETE Journal of Research, 1-17.

[24] Kadir, D. H. (2021). Statistical evaluation of main extraction parameters in twenty plant extracts for obtaining their optimum total phenolic content and its relation to antioxidant and antibacterial activities. Food Science & Nutrition, 9(7), 3491-3499.

[25] Wang, Z., Jin, Z., Yang, Z., Zhao, W., & Trik, M. (2023). Increasing efficiency for routing in Internet of Things using Binary Gray Wolf Optimization and fuzzy logic. Journal of King Saud University-Computer and Information Sciences, 35(9), 101732.

[26] Zhao, H., Zong, G., Zhao, X., Wang, H., Xu, N., & Zhao, N. (2023). Hierarchical Sliding-Mode Surface-Based Adaptive Critic Tracking Control for Nonlinear Multiplayer Zero-Sum Games Via Generalized Fuzzy Hyperbolic Models. IEEE Transactions on Fuzzy Systems.

[27] B. Moradhasel, A. Sheikhani, O. Aloosh, and N. J. Dabanloo, "Spectrogram classification of patient chin electromyography based on deep learning: A novel method for accurate diagnosis obstructive sleep apnea," Biomed Signal Process Control, vol. 79, p. 104215, 2023.

[28] Kadir, D. (2018). Bayesian inference of autoregressive models (Doctoral dissertation, University of Sheffield).

[29] Ding, X., Yao, R., & Khezri, E. (2023). An efficient algorithm for optimal route node sensing in smart tourism Urban traffic based on priority constraints. Wireless Networks, 1-18.

[30] Zhao, H., Zong, G., Wang, H., Zhao, X., & Xu, N. (2023). Zero-Sum Game-Based Hierarchical Sliding-Mode Fault-Tolerant Tracking Control for Interconnected Nonlinear Systems via Adaptive Critic Design. IEEE Transactions on Automation Science and Engineering.

[31] Fakhri, P. S., Asghari, O., Sarspy, S., Marand, M. B., Moshaver, P., & Trik, M. (2023). A fuzzy decision-making system for video tracking with multiple objects in non-stationary conditions. Heliyon.

[32] Zhang, H., Zou, Q., Ju, Y., Song, C., & Chen, D. (2022). Distance-based support vector machine to predict DNA N6-methyladenine modification. Current Bioinformatics, 17(5), 473-482.

[33] Saleh, D. M., Kadir, D. H., & Jamil, D. I. (2023). A Comparison between Some Penalized Methods for Estimating Parameters: Simulation Study. QALAAI ZANIST JOURNAL, 8(1), 1122-1134.

[34] Xue, B., Yang, Q., Jin, Y., Zhu, Q., Lan, J., Lin, Y., ... & Zhou, X. (2023). Genotoxicity Assessment of Haloacetaldehyde Disinfection Byproducts via a Simplified Yeast-Based Toxicogenomics Assay. Environmental Science & Technology, 57(44), 16823-16833.

[35] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1492–1500.

[36] A. Achille and S. Soatto, "Emergence of invariance and disentanglement in deep representations," The Journal of Machine Learning Research, vol. 19, no. 1, pp. 1947–1980, 2018.

[37] Pitchai, R., Supraja, P., Victoria, A. H., & Madhavi, M. J. N. P. L. (2021). Brain tumor segmentation using deep learning and fuzzy K-means clustering for magnetic resonance images. Neural Processing Letters, 53, 2519-2532.

[38] VK, D. (2023). An intelligent brain tumor segmentation using improved Deep Learning Model Based on Cascade Regression method. Multimedia Tools and Applications, 82(13), 20059-20078.

# Brightness Equalization Algorithm for Chinese Painting Pigments in Low-Light Environment Based on Region Division

Lijuan Cheng

College of Art and Art Design, Henan Vocational University of Science and Technology, Zhoukou, 466000, China

*Abstract*—With the promotion and development of Chinese painting and the advancement of photography technology, people can appreciate various types of Chinese paintings through image and other methods. However, Chinese painting images in low-light environments face the problem of extreme uneven brightness distribution. The currently proposed solutions for this problem are not sufficient. Therefore, this research proposes a brightness equalization algorithm for Chinese painting pigments in low-light environments based on region division. This algorithm also utilizes guided filtering for image denoising. In performance testing, the proposed method has a runtime of 16.63 seconds under a scaling factor of 1 and a runtime of 8.37 seconds under a scaling factor of 0.1, which are the fastest among the compared algorithms. In simulation experiments, the brightness equalization value of the proposed method is 198.93, which is listed at the best among all the compared algorithms. This research provides a valuable research direction for the brightness equalization of Chinese painting pigments.

*Keywords*—*Chinese painting; low-light; region division; guided filtering; scaling factor*

## I. INTRODUCTION

In the process of promoting Chinese painting, pigments are one of the most visually impactful elements for viewers, as they directly determine the visual effects of the artwork [1]. However, due to the limitations of low-light environments, Chinese paintings in low-light scenes may face the problem of uneven pigment brightness, which affects the visibility and artistic quality of the artwork to some extent. A low-light environment refers to an environment with dim lighting or insufficient light sources [2]. In such environments, due to the scarcity and weakening of light, details in the image are difficult to display clearly, and the differences in pigment brightness become more pronounced. This significantly affects the appreciation of Chinese paintings by viewers. At present, the processing methods for low light images include image enhancement, noise removal, and multi frame image fusion, which can be divided into two approaches: hardware and software. However, the focus is mainly on software upgrades and improvements [3]. Due to the uneven lighting in the shooting environment, the brightness of captured images may be uneven, with some areas appearing too bright while others are too dark. Existing brightness equalization algorithms not only enhance the noise in the image but also have issues with inconsistent equalization across different regions [4].

Therefore, this research proposes a brightness equalization

algorithm for Chinese painting pigments in low-light environments based on region division. This study will provide reference and guidance for the application of region-based brightness equalization algorithms in the field of Chinese painting pigment brightness equalization and other related fields. The research is divided into five sections: an overview of the research in Section I, a review of domestic and foreign studies in Section II, a study on the algorithm's methodology in Section III, performance testing of the algorithm in Section IV, and a summary and outlook on the limitations of this research in Section V.

## II. LITERATURE REVIEW

Researchers have focused on using image region division techniques to improve image enhancement methods. Matsuyama E proposed a segmentation method for chest X-ray images. This method can automatically remove the scapular region, mediastinal region, and diaphragm region from various chest X-ray images as the learning data for this method. The method uses a simple linear iterative clustering algorithm and local entropy filtering to generate an entropy map, which is then subjected to morphological operations to perform region segmentation on the lung image. The method was tested, and the results showed that it can remove non-pulmonary markings from the image and present clear X-ray images of the lungs [5]. Chen et al. found that the evaluation metrics of existing iris segmentation algorithms may be influenced by inaccurate localization of the Ground Truth image. Therefore, researchers proposed using a mask image segmented based on deep learning algorithms as a substitute for the Ground Truth image. Experimental results showed that the mask image segmented based on deep learning algorithms can completely replace the original Ground Truth image [6]. Cao et al. found that existing line art coloring methods can produce credible coloring results, but these methods are often affected by color bleeding issues. Therefore, researchers proposed an explicit segmentation fusion mechanism. Testing outcomes shows that the model can better fulfill the coloring instructions given by the user and can greatly alleviate the problem of color bleeding artifacts [7].

Image enhancement methods are applicable in various fields, and researchers have made many improvements to these methods. Tirumani et al. found that existing image enhancement methods have unstable effects on contrast and resolution enhancement. Therefore, researchers process the

resolution of the image, and then used auto optimization to enhance the resolution and brightness. Testing outcomes showed that this method can effectively and stably enhance the resolution and contrast of the image [8]. Xu et al. proposed a multi-scale fusion framework for low-light image enhancement. This framework first generates multiple artificially multi-exposure images using a mapping function, then combines exposure to create a weight map, and finally fuses different frequency bands of the image. Testing outcomes showed that this method outperforms existing algorithms in enhancing low-light images [9]. Lu et al. believed that the current image enhancement methods based on convolutional neural network models do not differentiate image features on different channels, which hinders the learning of hierarchical features. Therefore, researchers proposed a channel-split attention network that can analyze shallow features in a targeted manner by splitting them into residual and dense branches. Experimental results showed that this method exhibited excellent performance in both qualitative and quantitative evaluations [10].

In summary, although region division methods have been applied in multiple fields, their combination with image enhancement for brightness equalization of Chinese painting pigments is relatively rare. Therefore, this research proposes a brightness equalization algorithm for Chinese painting pigments in low-light environments based on region division, providing effective technical support for the promotion of Chinese painting.

## III. Region-based Brightness Equalization Algorithm for Chinese Painting Pigments in Low-Light Environments

Chinese painting images in low-light environments often suffer from uneven lighting in the pigment display [11]. Traditional image brightness equalization algorithms have limitations in achieving sufficient brightness equalization and enhancing all local details [12]. To address this issue, this research proposes a region-based brightness equalization algorithm for Chinese painting pigments in low-light environments. This algorithm improves existing image enhancement algorithms and provides effective assistance in the refinement of image enhancement algorithms.

### A. Single-Frame Image Brightness Equalization Enhancement Method in Low-Light Environments

Images in image processing come in various formats, including RGB, LAB, YUV, HSV, HIS, etc. The main image format studied in this research is HSV. In this research, the RGB image is first converted to the HSV image. In HSV, H represents the hue of the image, S represents the saturation, and V represents the value or brightness of the image. The advantage of the HSV format is that the color information of the image does not affect the brightness component, which ensures that the original colors of the image are preserved during brightness enhancement. The schematic diagram of an HSV image is shown in Fig. 1.

In Fig. 1, a cone shape is hired to manifest the HSV color space, where the hue is determined by the rotation angle around the center of the cone, with each 120° representing a

different color. The closer to the center of the cross-section, the less saturated the color is, and the closer to the apex of the cone, the weaker the brightness. The conversion of RGB to HSV is represented by Eq. (1).

$$V = \max\{R, G, B\} \tag{1}$$



Fig. 1. The schematic diagram of an HSV image.

In Eq. (1), $R, G, B$ represent the three primary colors in the RGB color space, where $R$ means red, $B$ represents blue, and $G$ represents green. The converted $S$ value is represented by Eq. (2).

$$S = \begin{cases} \dfrac{V - \min\{R, G, B\}}{V}, V \neq 0 \\ 0, otherwise \end{cases} \tag{2}$$

In Eq. (2), $V$ represents the value obtained from Eq. (1). The converted $H$ value is represented by Eq. (3).

$$H = \begin{cases} 0°, V = \min(R, G, B) \\ 60° \times \dfrac{G - B}{(V - \min(R, G, B))} + 0°, V = R, G \geq B \\ 60° \times \dfrac{G - B}{(V - \min(R, G, B))} + 360°, V = R, G < B \\ 60° \times \dfrac{G - B}{(V - \min(R, G, B))} + 120°, V = G \\ 60° \times \dfrac{G - B}{(V - \min(R, G, B))} + 240°, V = B \end{cases} \tag{3}$$

In Eq. (3), $R$ represents the $R$ value in the RGB. The brightness of an HSV image is divided into different levels [13]. The average brightness range of the V channel is $[0,1]$. An empirical threshold $I_{th}$ is set, and when the average brightness of an image is below this threshold, the image is

considered a high-brightness image. Therefore, the brightness range of high-brightness images is $[I_{th},1]$. For the enhancement of high-brightness images, the focus is mainly on enhancing the contrast [14]. Since the enhancement results of high-brightness images are similar to those of low-brightness images, this research converts high-brightness images to low-brightness images for enhancement and then converts them back to high-brightness images after enhancement [15]. The formula for obtaining the limited brightness image is represented by Eq. (4).

$$I_{\lim} = \begin{cases} I_v, \overline{I_v} \leq I_{th} \\ 1 - I_v, \overline{I_v} > I_{th} \end{cases} \qquad (4)$$

In Eq. (4), $I_{\lim}$ represents the limited brightness channel image, $I_v$ represents the $V$ channel image extracted after converting from RGB to HSV, and $\overline{I_v}$ represents the average brightness of the $V$ channel image. The brightness region division of the image in this research is divided into four steps. The first step is the initial enhancement of the limited $V$ channel image, and the enhancement formula is represented by Eq. (5).

$$F = \log_2(1 + I_{\lim}) \qquad (5)$$

In Eq. (5), $I_{\lim}$ represents the limited $V$ channel image, and $F$ represents the image after initial enhancement. The region segmentation of the image in this research is divided into four steps. The second step is the binarization of the multi-scale image [16]. Two different binarization methods are used for the edges and region shapes of the image. The first binarization method first applies mean filtering to the image $F$ after initial enhancement to obtain the neighborhood mean value of each pixel. Then, the brightness value is divided by the neighborhood mean value, and the result is compared with the adaptive sensitivity factor $T$. Finally, the binarized $V$ channel brightness image is obtained. The calculation process is represented by Eq. (6).

$$F_{binary\_1}(x, y) = \begin{cases} 1, \dfrac{F(x, y)}{F_{s1 \times s1}(x, y)} > T \\ 0, otherwise \end{cases} \qquad (6)$$

In Eq. (6), $F(x, y)$ represents the brightness value of each pixel, and $F_{s1 \times s1}(x, y)$ represents the neighborhood mean value. The second binarization method subtracts the mean filtering image from $F$, and then subtracts a constant $C$ to obtain the difference image $I_{sm}$. Then, based on the pixel values of I, binarization is performed to obtain a binary image containing texture boundaries. The calculation process is represented by Eq. (7).

$$F_{binary\_2}(x, y) = \begin{cases} 1, I_{sm}(x, y) > 0 \\ 0, otherwise \end{cases} \qquad (7)$$

In Eq. (7), $I_{sm}(x, y)$ represents the pixel values of the image $I_{sm}$. The second step of region image processing is fusion, which involves merging the two binarized images obtained earlier to create a new binary image. The fusion formula is shown in Eq. (8).

$$F_{binary} = F_{binary\_1} \oplus F_{binary\_2} \qquad (8)$$

In Eq. (8), $\oplus$ represents the logical AND operator. The third step is noise reduction using morphology. The fourth step is region segmentation. The specific operation involves first marking the boundaries of the denoised regions, and then dividing the image into multiple regions and assigning them numbers based on the marked content. The schematic diagram of image segmentation is shown in Fig. 2.

In Fig. 2, the images from left to right are the original image, the two binarized images, the fused image obtained from the fusion calculation of the two binarized images, the denoised binary image, and the segmented image. Due to the significant brightness differences between different regions in images with uneven lighting, targeted brightness adjustment is needed [17]. The image is marked based on the different brightness levels in different regions, and the marking rule is shown in Eq. (9).

$$\begin{cases} if \dfrac{(V_{i\min} + V_i)}{2} > 0.5, bright \\ else, dark \end{cases} \qquad (9)$$



Fig. 2. Schematic diagram of image partitioning.

In Eq. (9), $i$ represents the index of the region, $V_{i\min}$ represents the minimum brightness value in that region, and $\overline{V_i}$ represents the average brightness of that region.

### B. Image Brightness Equalization Enhancement Method Based on Region Denoising

Denoising is an important step in image enhancement. Currently, there are various denoising methods, including bilateral filtering-based denoising, Gaussian filtering-based denoising, and linear guided filtering-based denoising [18]. The denoising method based on bilateral filtering is computationally complex and slow. The denoising method based on Gaussian filtering tends to blur the edges of the denoised image and result in unclear presentation of image details. The denoising method based on linear guided filtering produces clear edges in the denoised image without artifacts and has a faster computation speed. Therefore, in this research, the guided filtering method is used for image denoising. The workflow of the guided filtering denoising method is shown in Fig. 3.



Fig. 3. Guiding the workflow of filtering and noise reduction methods.

In Fig. 3, the workflow of guided filtering includes a guidance image $I$, an input image $p$, and an output image $q$. The guidance image can be pre-set based on different application scenarios, but it can also be replaced by the input image. The guided filtering principle is based on the premise assumption that a linear relationship exists between the guidance image and the output image. Assuming that in a window $\omega_m$ centered at pixel $m$, $q$ is a linear transformation of $I$, the transformation formula is shown in Eq. (10).

$$q_i = a_m I_i + b_m, \forall i \in \omega_m \tag{10}$$

In Eq. (10), $a_m$ and $b_m$ represent the assumed linear invariant coefficients within the window $\omega_m$, and $\omega_m$ is a square window with a radius of $r$ centered at pixel $m$. To determine the values of $a_m$ and $b_m$, constraints need to be applied to the input image $p$, and the constraints are shown in Eq. (11).

$$q_i = p_i - n_i \tag{11}$$

In Eq. (11), $n$ represents the excess information in $q$,

where most of the irrelevant information is noise. The linear regression model established in window $\omega_m$ is shown in Eq. (12).

$$E(a_m, b_m) = \sum_{i \in \omega_m} ((a_m I_i + b_m - p_i)^2 + \varepsilon a_m^2) \tag{12}$$

In Eq. (12), $\varepsilon$ is a regularization parameter to constrain $a_m$, and its solution is shown in Eq. (13).

$$\begin{cases} a_m = \dfrac{\dfrac{1}{|\omega|}\sum_{i \in \omega_m} I_i p_i - \mu_m \overline{p}_m}{\sigma_m^2 + \varepsilon} \\ b_m = \overline{p}_m - a_m \mu_m \end{cases} \tag{13}$$

In Eq. (13), $\sigma_k^2$ and $\mu_k$ represent the variance and mean of $I$ within window $\omega_m$, $|\omega|$ represents the number of pixels in $\omega_m$, and $\overline{p}_m = \dfrac{1}{|\omega|}\sum_{i \in \omega_m} p_i$ represents the mean of $p$ within window $\omega_m$. Since pixel $i$ is included in different $\omega_m$, there will be multiple values of $q_i$ in different windows. Therefore, Eq. (14) is used to determine the value of $q_i$.

$$q_i = \frac{1}{|\omega|}\sum_{m|i \in \omega_m} (a_m I_i + b_m) = \overline{a}_m I_i + b_m \tag{14}$$

Due to the rotational symmetry property of the summation window, $\sum_{m|i \in \omega_m} a_m = \sum_{i|m \in \omega_i} a_m$ can be obtained. Therefore, Eq. (14) can also be written as Eq. (15).

$$q_i = \overline{a}_i I_i + \overline{b}_i \tag{15}$$

In Eq. (15), $\overline{a}_i = \dfrac{1}{|\omega|}\sum_{m \in \omega_i} a_m$ and $\overline{b}_i = \dfrac{1}{|\omega|}\sum_{m \in \omega_i} b_m$ represent the mean of the calculated results of the linear coefficients for pixel $i$. The denoising effects of different denoising methods on the same image are shown in Fig. 4.

In Fig. 4, Fig. 4(b) is the image after denoising with Gaussian filtering, Fig. 4(c) is the image after denoising with bilateral filtering, and Fig. 4(d) is the image after denoising with guided filtering. Fig. 4(a) is the original one. The image after denoising with Gaussian filtering appears darker in color and has unclear edges. The image after denoising with bilateral filtering has clear brightness edges. The image after denoising with guided filtering has clear brightness edges and the details in each region are smoothed, which better matches the actual lighting distribution [19]. In order to achieve better enhancement of the image, this research uses a two-dimensional gamma function to perform brightness correction on images with uneven lighting [20]. The formula for the two-dimensional gamma function is shown in Eq. (16).

$$I_g(x,y) = 255 \times \left( \frac{F(x,y)}{255} \right)^{\gamma}, \gamma = M(x,y)^{\frac{M(x,y)-I_q(x,y)}{M(x,y)}} \quad (16)$$



(a) Original drawing    (b) Gaussian filtered illumination map

(c) Bilateral filtered illumination map    (d) Guided filtered illumination map

Fig. 4.    Noise reduction effect diagram.

In Eq. (16), $F$ represents the preliminarily enhanced image after transformation, $I_g$ represents the output enhanced image, $M$ represents the target mean, and

$I_g(x,y)$ represents the image after denoising with guided filtering. When $M(x,y) > I_q(x,y)$ I, $\gamma < 1$, indicating an increase in brightness for $F(x,y)$; otherwise, the brightness is reduced. The image obtained is further filtered using guided filtering with a radius of six and a regularization parameter of $1e-4$ for denoising. Then, contrast-limited histogram equalization is applied to obtain the image $I_c$, which is then weighted fused using the weighted fusion formula shown in Eq. (17).

$$I_{out} = a \cdot I_c + b \cdot I_g, \begin{cases} a = \begin{cases} 0.5 + 0.1 \times \dfrac{\mu(I_v)-0.4}{0.6}, \mu(I_v) \geq 0.4 \\ 0.5, else \end{cases} \\ b = 1-a \end{cases} \quad (17)$$

In Eq. (17), $a$ represents the weighting coefficient for the output of contrast-limited histogram equalization, $b$ represents the weighting coefficient for $I_g$, and $\mu(I_v)$ represents the mean brightness of the $V$ channel image of the original image. Finally, the $I_{out}$ channel image, the $H$ (hue) channel image $I_H$ from the initial input image, and the $S$ (saturation) channel image $I_S$ are combined to form an HSV image, which is then converted to the RGB format and output as the final RGB image. In summary, the workflow of the brightness equalization algorithm based on region division in low-light conditions is shown in Fig. 5.



Fig. 5.    Flow chart of brightness equalization algorithm for Chinese painting pigments in low illumination environments based on region division.

In Fig. 5, the first step is to input the original image in RGB format. The second step is to convert the original RGB image to the HSV format and extract the V channel image $I_V$. The third step is to calculate the grayscale mean of image $I_V$ and determine if the mean is greater than a set empirical threshold $I_{th}$. If the mean is less than $I_{th}$, the mean remains unchanged. If the mean is greater than $I_{th}$, the mean is inverted, resulting in a mean-limited grayscale image $I_{\lim}$. The fourth step consists of three operations. The first operation is logarithmic transformation followed by binarization to obtain a binary image containing texture edges. Then, the binary image is segmented into regions. The second operation is 8-neighborhood mean filtering to obtain neighborhood information for each pixel. The third operation is denoising of the grayscale image $I_{\lim}$ using guided filtering, resulting in an illumination map $I_q$. The fifth step of the algorithm is to construct the target mean $M_i$ for each region based on the brightness mean, minimum value, and image mean of the original input image within the regions segmented in the first operation of the fourth step. The sixth step is to perform gamma correction on the preliminarily enhanced image to achieve brightness correction, resulting in the image $I_g$. The seventh step is to determine whether the inversion operation was performed in the third step. If inversion was performed, the image $I_g$ is restored; if not, it remains unchanged. The eighth step is to denoise the image $I_g$ using guided filtering to obtain the denoised image. The ninth step is to perform contrast-limited histogram equalization on $I_g$ and then perform weighted fusion to obtain the corrected image $I_{out}$. The final step is to combine the H and S channels back into the HSV color space, convert it to RGB format, and output the brightness equalized image.

## IV. PERFORMANCE ANALYSIS AND SIMULATION EXPERIMENT OF REGION-BASED IMAGE BRIGHTNESS EQUALIZATION ALGORITHM

### A. Performance Analysis of Region-based Image Brightness Equalization Algorithm

The processor used for this performance test is an Intel(R) Core (TM) i9-13900HX CPU with a clock speed of 5.4GHz, 16GB RAM, and a 64-bit operating system. The simulation software used is MATLAB R2022a. The images used in this experiment were obtained by continuously capturing 300 frames of the same Chinese painting in a low-light indoor environment. Five frames were randomly selected from the 300 frames of Chinese painting images and named Frame 1 to 5. The information entropy and Structural Similarity (SSIM) of the images enhanced by different algorithms were compared. The comparison of information entropy is Table I.

From Table I, the original images information entropy is generally in the range of 4-6. After MSRCR processing, it is in the range of 5-7. After Dong processing, the information entropy is in the range of 6-8. After Zohair processing, the information entropy is in the range of 7-9. After processing with the proposed method, the information entropy of the images is in the range of 9-10. A higher information entropy value indicates more detailed information in the image. The images processed by the proposed algorithm in this study show significantly more detailed information compared to other algorithms. The comparison results of SSIM are shown in Table II.

TABLE I. INFORMATION ENTROPY

| Algorithm | Frame 1 | Frame 2 | Frame 3 | Frame 4 | Frame 5 |
|---|---|---|---|---|---|
| Original drawing | 4.63 | 5.55 | 5.16 | 4.86 | 5.03 |
| MSRCR | 6.36 | 5.98 | 5.66 | 5.13 | 6.05 |
| Dong | 7.32 | 6.98 | 7.55 | 6.77 | 7.09 |
| Zohair | 7.53 | 8.25 | 7.86 | 8.03 | 7.33 |
| This study | 9.56 | 9.43 | 9.36 | 9.78 | 9.89 |

TABLE II. STRUCTURAL SIMILARITY COMPARISON RESULTS

| Algorithm | Frame 1 | Frame 2 | Frame 3 | Frame 4 | Frame 5 |
|---|---|---|---|---|---|
| MSRCR | 0.06 | 0.13 | 0.22 | 0.09 | 0.23 |
| Dong | 0.23 | 0.15 | 0.09 | 0.21 | 0.26 |
| Zohair | 0.22 | 0.26 | 0.19 | 0.23 | 0.25 |
| This study | 0.30 | 0.34 | 0.29 | 0.36 | 0.35 |

From Table II, the SSIM of the images reinforced by MSRCR is in the range of 0.06-0.23. The SSIM of the images enhanced by Dong is in the range of 0.09-0.26. The SSIM of the images enhanced by Zohair is in the range of 0.19-0.26. The SSIM of the images enhanced by the proposed method in this study is in the range of 0.30-0.36. The SSIM of the images enhanced by the method is significantly higher than other algorithms, indicating that the details of the images preserved by the method are more complete and the enhancement effect is better compared to other algorithms. A comparison was made between the adaptive gamma brightness correction method used in the algorithm and the fixed parameter gamma correction method. The experimental results are shown in Fig. 6.



Fig. 6. Comparison of different gamma brightness correction methods.

In Fig. 6, from the experimental results, it can be observed that when the input images are in the range of 1-300 frames, the brightness fluctuation of the images corrected by the

improved adaptive gamma correction method in this study significantly decreases compared to the original images. On the other hand, the images corrected by the fixed parameter gamma correction method exhibit larger brightness fluctuations compared to the original images. The adaptive gamma correction method used in this study outputs a standard deviation of brightness of $1.861\times10^{-3}$, while the original images have a brightness standard deviation of

$7.634\times10^{-3}$, and the images corrected by the fixed parameter gamma correction method have a brightness standard deviation of $5.342\times10^{-3}$. The standard deviation of brightness corrected by the adaptive gamma correction method is significantly lower than that of the original images and the fixed parameter gamma correction method. This study also compared the runtime of different algorithms at different scaling factors. The specific results are shown in Fig. 7.



Fig. 7. Running time of different algorithms under different scaling coefficients.

Fig. 7(a) represents the runtime of different algorithms with varying group numbers under a scaling factor of 1. It can be observed that under a scaling factor of 1, the runtime of all algorithms increases with the increase in group numbers. When the group number reaches 200, the proposed algorithm in this study has the shortest runtime among all algorithms, which is 16.63 seconds. Fig. 7(b) represents the runtime of different algorithms with varying group numbers under a scaling factor of 0.1. It can be seen that under a scaling factor of 0.1, the runtime of all algorithms is significantly shorter compared to the case of a scaling factor of 1. Among them, the proposed algorithm in this study has the shortest runtime among all group numbers, which is 8.37 seconds when the group number reaches 200.

### B. Simulation Experiment of Brightness Equalization Algorithm for Chinese Painting Images Based on Regional Division

The comparison of the average brightness of Chinese painting pigments enhanced by different enhancement algorithms in low-light environments is shown in Fig. 8.

In Fig. 8, the image enhanced by the MSRCR image enhancement algorithm has the lowest average brightness among the five algorithms, indicating that its brightness equalization processing is the worst among the five algorithms. The image enhanced by the Dong image enhancement algorithm has a lower average brightness and a slower growth rate. The image enhanced by the Zohair image enhancement algorithm has a higher average brightness and a faster growth rate. The image enhanced by the algorithm proposed maintains a high level of average brightness and has a fast growth rate. This study introduced the Peak Signal-to-Noise Ratio (PSNR) for evaluating the fidelity of the images. PSNR tests were

conducted on different algorithms using the Brightening dataset and the LOL dataset, and the results are shown in Fig. 9.



Fig. 8. Average brightness results of Chinese painting pigments in low illuminance environments with different enhancement algorithms for image enhancement.

In Fig. 9, Fig. 9(a) shows the comparison of PSNR for different algorithms under the Brightening dataset. The PSNR of the MSRCR algorithm is close to that of the Dong algorithm, and both algorithms have large fluctuations. The PSNR of the Zohair algorithm is higher, and its fluctuations are more stable than the above two algorithms. The PSNR of the algorithm proposed is greater than the above three at bit rates ranging from 0 to 2000, and it has a fast growth rate. The maximum PSNR values for the four algorithms are obtained when the bit rate reaches 2000, which are 26.6db, 28.4db,

32.5db, and 37.6db, respectively. Fig. 9(b) shows the comparison of PSNR for different algorithms under the LOL dataset. Except for the method proposed, the PSNR of the other three algorithms fluctuates significantly under the LOL dataset. However, the maximum PSNR values for each algorithm are still obtained at a bit rate of 2000, which are 28.2db, 32.6db, 35.8db, and 37.1db, respectively. The maximum PSNR for the four algorithms has improved under the LOL dataset. A comparison was made between the brightness histograms of the low-light Chinese painting images enhanced by different algorithms and the brightness histogram of the original Chinese painting image, and the results are shown in Fig. 10.

In Fig. 10, Fig. 10(a) represents the brightness distribution histogram of the original low-light Chinese painting image. It

can be seen that the high brightness region of the original image is concentrated in one area, and the brightness in other areas is at a lower level, indicating a highly uneven brightness distribution of the image. Fig. 10(b) represents the brightness distribution histogram after brightness equalization processing using the Zohair algorithm. It can be seen that the overall brightness of the image has significantly improved, but there are still some areas with low brightness values, indicating an uneven brightness distribution. Fig. 10(c) represents the brightness distribution histogram after brightness equalization processing using the algorithm proposed. The overall brightness of the image has significantly improved, and the brightness distribution is balanced in all parts, indicating that the proposed method can better perform brightness equalization processing on Chinese paintings captured in low-light environments.


(a) Brightening dataset


(b) LOL dataset

Fig. 9. Comparison of PSNR under different datasets.


(a) Histogram of the original image


(b) Histogram of the Zohair image


(c) The histogram of the algorithm proposed in this study

Fig. 10. Brightness histograms of low illuminance Chinese painting images enhanced by different algorithms.

## V. CONCLUSION

The uneven display of pigment brightness in low light environments in Chinese painting poses certain obstacles to the promotion of Chinese painting and the dissemination of Chinese culture. Image denoising and image enhancement are common methods for processing low light images. However, traditional image denoising methods such as mean filtering and bilateral filtering may lead to loss of image details and poor generalization performance [21]. Traditional image enhancement methods based on grayscale transformation, histogram equalization, and Retinex theory also suffer from poor visual effects, unsatisfactory enhancement effects, and loss of details [22-23]. In this context, in order to enhance low light images more effectively, a region-based brightness equalization algorithm for low-light Chinese paintings was proposed. The performance test tells that the image information entropy enhanced by the proposed method was between 9 and 10, higher than the 5-7 of MSRCR, the 6-8 of Dong, and the 7-9 of Zohair, reaching the highest value among all the compared algorithms. This indicates that the image processed by the proposed algorithm presents more detailed information compared to other algorithms. The SSIM (Structural Similarity Index) of the image enhanced by the proposed method was between 0.30 and 0.36, significantly higher than the 0.06-0.23 of MSRCR, the 0.09-0.26 of Dong, and the 0.19-0.26 of Zohair. This indicates that the details of the image preserved after enhancement using the proposed method are more complete compared to other algorithms, resulting in better enhancement effects. In the simulation experiment, the brightness distribution histograms of the images after brightness equalization processing using different algorithms were compared. The testing outcomes tells that overall brightness value of the image processed by the proposed method was significantly improved compared to the original image. The brightness values were roughly between 1.51 and 2.0, indicating a more balanced distribution. This indicates that the proposed method can effectively perform brightness equalization processing for Chinese paintings in low-light environments. However, the proposed brightness equalization algorithm is easily affected by image brightness and noise when performing region partitioning, and it is implemented through simulation on the Matlab platform, which limits its runtime. Therefore, further research can explore better denoising methods for image preprocessing, algorithm optimization, and GPU acceleration.

## REFERENCES

[1] Mishro P K, Agrawal S, Panda R, Abraham A. A novel brightness preserving joint histogram equalization technique for contrast enhancement of brain MR images. Biocybernetics and Biomedical Engineering, 2021, 41(2):540-553.

[2] Tirumani V H L, Tenneti M, Srikavya K C, Kotamraju S K. Image resolution and contrast enhancement with optimal brightness compensation using wavelet transforms and particle swarm optimization.IET image processing, 2021, 15(12):2833-2840.

[3] Jeevan K M, Anne G A B, Kumar P V. An image enhancement method based on gabor filtering in wavelet domain and adaptive histogram equalization. Indonesian Journal of Electrical Engineering and Computer Science, 2021, 21(1):146-153.

[4] Rao G S, Srikrishna A. Image Pixel Contrast Enhancement Using Enhanced Multi Histogram Equalization Method. Ingénierie des Systèmes D Information, 2021, 26(1):95-101.

[5] Hussain I, Muhammad J. Efficient convex region-based segmentation for noising and inhomogeneous patterns. Inverse Problems and Imaging, 2023, 17(3):708-725.

[6] Matsuyama E. A Novel Method for Automated Lung Region Segmentation in Chest X-Ray Images. Journal of biomedical science and engineering, 2021, 14(6):288-299.

[7] Chen Y, Gan H, Zeng Z, Chen H. DADCNet: Dual attention densely connected network for more accurate real iris region segmentation. International Journal of Intelligent Systems, 2021, 37(1):829-858.

[8] Cao R, Mo H, Gao C. Line Art Colorization Based on Explicit Region Segmentation. John Wiley & Sons, Ltd, 2021, 40(7):1-10.

[9] Tirumani V H L, Tenneti M, Srikavya K C, Kotamraju S K. Image resolution and contrast enhancement with optimal brightness compensation using wavelet transforms and particle swarm optimization.IET image processing, 2021, 15(12):2833-2840.

[10] Xu Y, Yang C, Sun B, Yan X, Chen M. A novel multi-scale fusion framework for detail-preserving low-light image enhancement. Information Sciences, 2021, 548(12):378-397.

[11] Liang Y. Analysis of the Integration of Chinese Painting Techniques in Watercolor Painting. Arts Studies and Criticism, 2022, 3(1):37-40.

[12] Lu B, Pang Z, Gu Y, Zheng Y. Channel splitting attention network for low-light image enhancement. IET Image Processing, 2022, 16(5):1403-1414.

[13] Reddy S K, Prasad R K. An Extended Fuzzy C-Means Segmentation for an Efficient BTD With the Region of Interest of SCP. International journal of information technology project management, 2021, 12(4):11-24.

[14] Wang W, Liu R. A saturation-value histogram equalization model for color image enhancement. Inverse Problems and Imaging, 2023, 17(4):746-766.

[15] Yu N, Li J, Hua Z. LBP-based progressive feature aggregation network for low-light image enhancement.IET image processing, 2022, 16(2):535-553.

[16] Choudhuri S, Adeniye S, Sen A. Distribution Alignment Using Complement Entropy Objective and Adaptive Consensus-Based Label Refinement for Partial Domain Adaptation. Artificial Intelligence and Applications. 2023, 1(1): 43-51.

[17] Yang Y, Song X. Research on face intelligent perception technology integrating deep learning under different illumination intensities. Journal of Computational and Cognitive Engineering, 2022, 1(1): 32-36.

[18] Ponmani E, Saravanan P. Image denoising and despeckling methods for SAR images to improve image enhancement performance: a survey. Multimedia Tools and Applications, 2021, 80(17): 26547-26569.

[19] Xu Y, Yang C, Sun B, Yan X, Chen M. A novel multi-scale fusion framework for detail-preserving low-light image enhancement. Information Sciences, 2021, 548(12):378-397.

[20] Sandoub G, Atta R, Ali H A, Abdel-KaderA R F. low-light image enhancement method based on bright channel prior and maximum colour channel. IET Image Processing, 2021, 15(8):1759-1772.

[21] Ilesanmi A E, Ilesanmi T O. Methods for image denoising using convolutional neural network: a review. Complex & Intelligent Systems, 2021, 7(5): 2179-2198.

[22] Hosny K M, Darwish M M, Aboelenen T. Novel fractional-order polar harmonic transforms for gray-scale and color image analysis. Journal of the Franklin Institute, 2020, 357(4): 2533-2560.

[23] Bulut F. Low dynamic range histogram equalization (LDR-HE) via quantized Haar wavelet transform. The Visual Computer, 2022, 38(6): 2239-2255.

# Anomaly Detection in Structural Health Monitoring with Ensemble Learning and Reinforcement Learning

Nan Huang*

Jiangsu University of Science and Technology, Zhenjiang 212000, Jiangsu, China

*Abstract*—**This research introduces a novel approach for improving the analysis of Structural Health Monitoring (SHM) data in civil engineering. SHM data, essential for assessing the integrity of infrastructures like bridges, often contains inaccuracies because of sensor errors, environmental factors, and transmission glitches. These inaccuracies can severely hinder identifying structural patterns, detecting damages, and evaluating overall conditions. Our method combines advanced techniques from machine learning, including dilated convolutional neural networks (CNNs), an enhanced differential equation (DE) model, and reinforcement learning (RL), to effectively identify and filter out these irregularities in SHM data. At the heart of our approach lies the use of CNNs, which extract key features from the SHM data. These features are then processed to classify the data accurately. We address the challenge of imbalanced datasets, common in SHM, through a RL-driven method that treats the training procedure as a sequence of choices, with the network learning to distinguish between less and more common data patterns. To further refine our method, we integrate a novel mutation operator within the DE framework. This operator identifies key clusters in the data, guiding the backpropagation process for more effective learning. Our approach was rigorously tested on a dataset from a large cable-stayed bridge in China, provided by the IPC-SHM community. The results of our experiments highlight the effectiveness of our approach, demonstrating an Accuracy of 0.8601 and an F-measure of 0.8540, outperforming other methods compared in our study. This underscores the potential of our method in enhancing the accuracy and reliability of SHM data analysis in civil infrastructure monitoring.**

*Keywords—Structural health monitoring; Anomaly detection; reinforcement learning; differential equation; imbalanced classification*

## I. INTRODUCTION

SHM is a key method for overseeing civil infrastructure, offering insights into structural loads, performance, responses, and future behavior predictions. SHM's widespread adoption has led to a significant increase in data generation; for example, China's Sutong Bridge, with its 785 sensors, produces 2.5 TB of data annually [1, 2]. Analyzing this vast amount of SHM data is challenging due to various anomalies caused by sensor errors, system failures, environmental factors, and more. These issues, compounded by data from significant events like earthquakes or accidents, can jeopardize the accuracy of structural analysis and the predictive power of SHM systems [3].

Implementing sensor-driven SHM methods generates large amounts of sequential data, complicating manual analysis and anomaly detection. Variations in this data may result from diverse factors such as weather, vehicle overloads, accidents, or unexpected events. It is crucial to recognize that not all anomalies indicate structural issues; some may stem from sensor errors, calibration issues, noise, or transmission problems. To tackle these anomalies, solutions can be applied at both hardware and software levels. While hardware solutions like using wired data channels, extra sensors, or self-validating sensors are effective, they are often costly. Therefore, there is a growing preference for advanced data preprocessing techniques specifically designed for anomaly detection.

SHM faces the challenge of data imbalance, where class instances vary significantly in number. To tackle this issue, two approaches are used: data-centric and algorithm-based. Data-centric strategies, such as under-sampling, over-sampling, and hybrid methods, aim to balance class distribution. Notably, the synthetic minority oversampling technique (SMOTE) [4] creates new minority class instances by linear interpolation, while NearMiss [5] under-samples the majority class using a nearest neighbor algorithm. However, over-sampling can lead to overfitting, and under-sampling may lose critical information. Algorithmic approaches focus on emphasizing the underrepresented class. These include modifying ensemble learning, altering decision thresholds, and employing cost-sensitive learning strategies. Cost-sensitive methods treat classification as cost minimization, assigning higher misclassification costs to minority cases. Ensemble methods combine multiple classifiers for a final decision, and threshold adjustment methods tweak the decision threshold during testing. These techniques aim to effectively balance accuracy and information retention in SHM data classification [6].

Furthermore, the incorporation of deep learning methodologies can serve as an avenue to address the challenge of imbalanced classification [7, 8]. Deep Reinforcement Learning (DRL) emerges as a promising solution to handle imbalanced data due to its distinct attributes. By employing a reward mechanism, DRL can assign augmented importance to the minority class, either by imposing stricter penalties for misclassifying instances from the minority class or by offering greater rewards for accurately identifying them. This approach actively counters the bias that conventional techniques display towards the majority class. The advantages of DRL extend beyond the mere balancing of class distribution. It also enriches the visibility of crucial patterns, particularly those associated with the minority class, by effectively filtering out noisy data. DRL's ability to unearth significant yet often overlooked features within the data contributes to the development of a more robust and efficient model [9].

The initial weight configuration in neural networks is crucial for training in SHM prediction. Traditional training, often using gradient-driven algorithms like backpropagation, typically starts with randomly assigned weights. However, the initial weight selection greatly impacts the training's efficiency and outcome. Careful consideration of the initialization strategy is essential for effective training and accurate SHM prediction. One effective approach is population-based training, where the best solution from a range of generated models is chosen as the starting point for the neural network. This method helps avoid the common issue of getting stuck in local optima, prevalent in standard training methods. Notably, simple evolutionary algorithms have shown effectiveness on par with stochastic gradient descent in neural network training [10].

DE [11] is a popular population-based optimization algorithm widely used in solving optimization problems, particularly effective for weight initialization in machine learning. DE offers several advantages: it ensures a broad exploration of the solution space, preventing entrapment in local optima and leading to better weight configurations. It updates weights iteratively based on the difference between current and target solutions, promoting faster convergence and improved performance. DE is also resilient to noise in fitness assessments, adeptly handling data uncertainties during weight initialization and providing stable initial weight settings. Furthermore, DE's flexibility and adjustability in weight initialization permit tailoring to particular problem areas, like establishing weight limits or integrating previous insights. This versatility improves DE's capability to initiate weights that are aptly matched to the distinct learning challenges being addressed.

This study investigates a novel approach combining a RL-based training algorithm with an advanced DE technique, specifically designed for SHM of bridges. It focuses on detecting anomalies in time-series sensor data from a major cable-stayed bridge in China. The data is divided into seven categories: normal, trend, square, missing, minor, drift, and outlier, with 'normal' being the most frequent. To overcome the issue of class imbalance in the dataset, the research introduces a framework that treats classification as a series of strategic decisions. In each iteration, an agent assesses a training sample (environmental state) and makes classification decisions, earning rewards or penalties based on the outcome. Classes with fewer samples are assigned higher rewards, encouraging accurate identification of less common anomalies. Additionally, the study integrates a unique mutation operator based on clustering principles within the DE framework to improve the backpropagation (BP) process. This operator identifies dominant clusters in the DE population and implements a novel approach for creating potential solutions. The key contributions of this research lie in its innovative approach to class imbalance, decision-making process in classification, and enhanced training methodology through integrating RL, DE, and BP processes:

*1) We* present an innovative RL-based method specifically designed to address the inherent challenges of imbalanced classification in SHM.

*2) The* approach integrates a unique reward system that reinforces accurate decisions while penalizing incorrect ones. By allocating enhanced rewards to the less represented class, we directly address the challenge of data skewness, encouraging the algorithm to appropriately focus on lesser-known data. This strategic maneuver contributes to a more fair and balanced classification procedure.

*3) To* extract deeper insights from images and refine the classification decision-making process, we employ a fusion of CNN models. This approach enhances the representation of features, resulting in improved accuracy and robustness in classification efforts.

*4) We* have developed an enhanced DE algorithm to initialize weights in the proposed model efficiently. This tactic aids in identifying a promising region for initiating the BP algorithm within the model.

The structure of this document is as follows: Section II details a review of relevant literature, and Section III provides an overview of the key dataset utilized in this study. Section IV delves into the proposed strategy, elucidating the core methodology in depth. Section V unfolds the empirical outcomes and their subsequent dissection. Concluding observations and potential avenues for future inquiry are encapsulated in Section VI.

## II. LITERATURE REVIEW

Artificial Intelligence (AI) techniques bring forth the capability to uncover patterns within time-series data with no prior comprehension of the underlying structural architecture. These methodologies involve exploring either the time or frequency domain of the data, extracting pertinent characteristics through statistical evaluations, or employing signal processing tools like the Fourier and wavelet transforms, as well as the Hilbert-Huang and Shapelet transforms. On the other hand, deep learning (DL) algorithms possess the ability to autonomously extract significant attributes by interpreting time-frequency data as visual inputs within a CNN framework. However, it is important to acknowledge that DL-centric approaches, while potent, demand substantial computational resources and cause meticulous fine-tuning of hyperparameters [12].

In order to tackle irregularities within SHM data, Pan et al. [13] presented an approach rooted in transfer learning. They employed a deep neural network to discern and rectify aberrant data, enhancing the accuracy of bridge evaluations. Samudra et al. [12] devised a comprehensible framework rooted in decision trees, employing random forest classifiers to categorize acceleration data in the realm of SHM. This approach, boasting a remarkable 98% accuracy, emerges as an economically viable avenue for gauging infrastructure state. Li et al. [14] outlined a strategy to elevate the efficacy of anomaly detection within bridge SHM systems. Employing strategies like data augmentation, feature dimension reduction, and a two-stage deep convolutional neural network, they achieved an elevated level of recognition accuracy. Tang et al. [3] presented an innovative anomaly detection method catering to SHM, employing a CNN that transforms time series data into visual representations. This approach achieves precise

identification of diverse pattern anomalies, scaling effectively and bolstering accuracy. Ye et al. [15] proposed a technique rooted in deep learning for identifying data anomalies within SHM systems. By deploying time-frequency analysis and CNNs, they translated SHM data into RGB images, subsequently classified through a Google network. Green et al. [16] introduced a new way to use Bayesian techniques in the analysis of inclinometer data for SHM. This method allows for the detection of anomalies, forecasting, and quantifying uncertainties, leading to better risk assessment and cost reduction.

Moreover, the framework has the potential to be applied in various engineering fields beyond inclinometers. Boccagna et al. [17] suggested an AI approach for monitoring structural health in almost real-time, using unsupervised deep learning. By preprocessing data and utilizing artificial neural network autoencoders, the technique effectively identifies anomalies, surpassing current methods and demonstrating encouraging outcomes. Lei et al. [18] proposed a residual attention network (RAN) to detect abnormal data in measured structures. The RAN incorporates attention mechanisms and residual learning to enhance classification accuracy and efficiency. It achieved exceptional performance and generalization on datasets from an arch bridge and a cable-stayed bridge, surpassing existing models in terms of multi-classification and accuracy. Yang and Nagarajaiah [19] introduced a principled, independent component analysis approach to reduce faulty data during data transmission; then, they achieved reliable data transfer and image restoration using compression sensing methods [20]. Yang and Nagarajaiah [21] employed the principal component pursuit method to detect and minimize burst noise in ambient vibration response. They also introduced a data management and processing framework based on sparsity rank and low-rank techniques. Park et al. [22] utilized transmission errors and ensemble empirical mode decomposition to identify anomalies such as gear teeth spalls and cracks in rotating machinery.

## III. DATASET DESCRIPTION

In this study, our primary focus centers on the detection of specific deviations - including trends, squares, omissions, minor variances, drifts, and outlier - within the acceleration time series derived from a lengthy cable-stayed bridge in China. The IPC-SHM community [23] provides access to curated data from this bridge. An overview of these anomalies, distilled from the bridge's measured data, is concisely summarized in Table I.

Besides the irregularities mentioned, it is vital to acknowledge that acceleration sensors, particularly those affixed to structures with potential vulnerabilities, adeptly detect a broad range of atypical patterns. These include offsets, characterized by sudden, noticeable jumps in response, and gains, marked by a slow, consistent increase in response over time. Furthermore, the sensors can identify precision deterioration, where the response shows erratic fluctuations, and complete failures, which result in a response akin to the randomness of white noise in the frequency domain. Recognizing and interpreting these additional types of deviations are crucial for comprehensive structural health monitoring, as they can provide early warning signs of more significant issues or impending failures.

The dataset under study contains a thorough record of acceleration data over a month, meticulously gathered from 38 strategically placed accelerometers across the bridge. For detailed analysis, this data is segmented into individual hourly time series, leading to an extensive collection of 28,272 such series. This figure is calculated considering the number of sensors, the days in the month, and the daily time cycle. Given that the accelerometers recorded at a rate of 20 Hz, the total data volume reaches an astonishing $2 \times 10^9$ data points, a multiplication of the number of time series, seconds in an hour, and the sampling rate. This vast dataset offers a rich source for in-depth analysis, allowing for the examination of minute changes and patterns over time, providing a comprehensive understanding of the bridge's dynamic behavior under various conditions.

The subsequent classification task systematically organizes these 28,272 time series responses into seven distinct categories. This includes the 'normal' set and six other types of anomalies: trend, square, missing, minor, drift, and outlier. A detailed chart in Table I itemizes the precise distribution of time series across each category within this dataset. This categorization is crucial for identifying the predominant anomalies and understanding their relative occurrences. It aids in developing targeted strategies for monitoring and maintenance, ensuring focused attention on the most critical or frequently occurring issues, enhancing the overall efficiency and effectiveness of the structural health monitoring process.

TABLE I. DESCRIPTION OF THE ANOMALIES

| Class | Description | Count |
|---|---|---|
| Normal | The time-domain response showcases balance, while multiple resonance peaks can be observed in the frequency-domain response. | 13575 |
| Trend | A discernible trajectory is apparent in the time-domain response, and a distinctive peak value is identified in the frequency-domain response. | 5778 |
| Square | The time-domain response mirrors a square wave. | 2996 |
| Missing | The bulk of the time-domain response is absent. | 2942 |
| Minor | Compared to standard category data, the time-domain response exhibits a notably reduced magnitude. | 1775 |
| Drift | The time-domain response varies unpredictably, either exhibiting arbitrary shifts or increasingly diverging over time. | 679 |
| Outlier | The time-domain response contains singular or multiple pronounced protrusions. | 527 |

## IV. MODEL ARCHITECTURE

Fig. 1 illustrates the intricacies of our innovatively developed model, meticulously engineered to boost the efficiency of anomaly detection. This model is specifically tailored to tackle issues like uneven distribution of classes and the critical importance of accurately setting initial weights. By integrating the DE algorithm with PPO, our model adeptly circumvents common hurdles encountered in standard models.

Conventional methods often miss a systematic strategy for establishing initial weights, resulting in reduced learning speeds and a propensity to settle on less-than-ideal outcomes. Our technique leverages DE to provide a diverse spectrum of initial weights. This variety helps the model to bypass smaller optima and more efficiently converge on more comprehensive solutions. Expanding on this, the utilization of DE in our method is not just about broadening the range of initial weights, but also about introducing a dynamic and adaptive way of initializing these weights. This adaptability is crucial in complex models where the landscape of solutions is vast and varied. By starting from a more helpful position in the solution space, our method enhances the model's ability to navigate through this landscape, leading to quicker and more effective convergence. Furthermore, this approach also contributes to the robustness of the model, making it less susceptible to the challenges posed by different data distributions and complexities inherent in various learning tasks. Our method does not just improve the efficiency of the learning process, but also broadens its applicability and effectiveness across a range of scenarios.

Additionally, the RL element in our framework is thoughtfully structured to significantly favor the accurate identification of the minority class, underscoring these crucial predictions. This represents a significant advancement over conventional supervised learning approaches, which often struggle with insufficient data representation across various classes. The dynamic nature of policy learning in RL encourages a fairer approach to decision-making, leading to enhanced strategies for identifying underrepresented categories. The flexibility of RL in our setup sets it apart from orthodox methodologies, equipping it with the essential capabilities to efficiently address the typical challenges found in conventional classification techniques employed in anomaly detection.



Fig. 1. Our model pipeline encompasses a sequence of procedures.

### A. Pre-training Phase

Deep models depend heavily on the initialization of deep network weights. If the initial values are not accurate, it can lead to convergence issues in the model. The first stage in this paper is to set the CNN and feed-forward neural network weights. We offer a more effective DE approach, which incorporates the power of a clustering technique and an innovative fitness function to optimize its performance. In our improved DE algorithm, we utilize a mutation and updating scheme based on clustering to enhance the optimization performance. Complex architectures rely significantly on the initial setup of deep network parameters. Inaccurate starting

points may cause the model to struggle with convergence. The initial step in our study involves configuring the weights for the CNN and the feed-forward neural networks. We propose an enhanced DE method that integrates the effectiveness of a clustering algorithm and a novel fitness function to boost its efficiency. In our refined DE technique, we employ a mutation and refresh strategy centered on clustering to improve the optimization process. Extending this, our approach takes into account the intricacies of deep learning architectures, ensuring that the weight initialization is not only randomly, but strategically influenced by the underlying data structure. The clustering-based mutation strategy allows for a more targeted and data-driven adjustment of weights, which is particularly useful in navigating the high-dimensional spaces typical in deep learning. This method aids in avoiding local minima and accelerates the convergence of the model.

The mutation mechanism, influenced by studies referenced in [24], identifies a promising region in the exploration domain. Employing the k-means clustering technique, the existing group P is segregated into k segments, each representing a unique portion of the exploration zone. A random integer chosen from the interval $[2, \sqrt{N}]$ dictates the cluster count. The cluster deemed most optimal possesses the lowest mean fitness value across its gathered samples post-clustering. Expanding on this, our method enhances the search strategy within the algorithm. By dividing the population into clusters, we can pinpoint specific regions in the search space that hold potential for better solutions. This clustering not only focuses the search but also adds a layer of precision in identifying promising areas, thereby increasing the efficiency of the mutation process. Additionally, the number of clusters is dynamically determined based on the population size, allowing for a flexible and adaptable approach to clustering. This adaptability is crucial in dealing with diverse problems and varying sizes of search spaces. The concept of assessing the perfection of a cluster based on its average fitness introduces a competitive element among the clusters, driving the algorithm to favor areas of the search space that show higher potential for optimal solutions. Furthermore, our approach refines the selection process within each cluster. After identifying the most optimal cluster, we focus on fine-tuning the solutions within this cluster, leveraging the collective intelligence of the group. This targeted mutation within the most promising cluster ensures the algorithm does not just wander aimlessly across the entire search space but makes informed, strategic moves towards areas more likely to yield superior results.

The clustering-based approach outlines the proposed mutation:

$$\overrightarrow{v_i^{clu}} = \overrightarrow{win_g} + F(\overrightarrow{x_{r_1}} - \overrightarrow{x_{r_2}}) \tag{1}$$

where, $\overrightarrow{x_{r_1}}$ and $\overrightarrow{x_{r_2}}$ represent two candidate solutions randomly selected from the current population while $\overrightarrow{win_g}$ corresponds to the best solution within the promising region. It is important to note that $\overrightarrow{win_g}$ may not be the best solution for the entire population.

Following the creation of M new solutions via mutation grounded in clustering, the current population undergoes an update under GPBA [25]. The GPBA is an innovative optimization approach that meticulously navigates through the search space, assessing the effectiveness of solutions against a series of established patterns. In practical application, GPBA changes the population by meticulously choosing and substituting individuals according to their performance indicators. By employing these patterns as navigational aids, GPBA refines its quest for the best solutions, often resulting in a more streamlined and efficient path to the global optimum. This method's real strength lies in its structured approach to the exploration of the search space. By using gradient patterns, the algorithm can intelligently predict the direction in which improvements can be made, rather than relying on random or exhaustive search methods. This predictive capability is helpful in complex optimization scenarios, characterized by vast search spaces and elusive optimal solutions. It allows the algorithm to bypass fewer promising regions of the search space, focusing its efforts on areas more likely to yield fruitful results. Furthermore, the GPBA's adaptability to different optimization problems adds to its versatility. Whether the task involves continuous or discrete variables, linear or non-linear relationships, the GPBA can be tailored to suit the specific characteristics of the problem. This adaptability is achieved through the customization of its pattern-based search mechanisms, which can suit various problem structures and complexities. Besides its efficiency in finding solutions, the GPBA also offers improved computational speed compared to more traditional optimization methods. This is beneficial in real-time applications or scenarios where time is a critical factor. The algorithm's ability to quickly converge to an optimal solution without sacrificing accuracy makes it an attractive choice for a wide range of optimization tasks.

The process unfolds as follows:

- Selection: To initiate the algorithm, generate k random individuals that will function as the initial points.

- Generation: Produce a set of $M$ solutions using mutation based on clustering and denote it as $v^{clu}$.

- Replacement: Choose $M$ solutions randomly from the current population to form set $B$.

- Update: Select the top M solutions from the combined groups $v^{clu}$ and B to create a new group $B'$. The refreshed population is derived by merging members of set P not included in B with those from set $B'$ $((P - B) \cup B')$.

*B. DRL*

DRL stands as a formidable approach in the domain of deep learning. Within this framework, an intelligent agent engages dynamically with its environment, aiming to maximize its cumulative rewards. This flexible and adaptive learning mechanism empowers the agent to make a series of decisions, often in the face of uncertainty, which has profound applications across a wide spectrum of domains, including but not limited to robotics, healthcare, and finance [26]. The prowess of DRL becomes evident in tasks that require sequential decision-making and the ability to adapt to unforeseen and evolving circumstances. Its capacity to handle complex activities that unfold over time, adjusting its strategies

and responses as needed, underscores its versatility and broad applicability in addressing real-world challenges. DRL's ability to learn from interactions with the environment, optimize decision-making processes, and navigate through dynamic scenarios positions it as a valuable tool for a wide range of applications, making it a compelling area of research and development in artificial intelligence.

In categorization-related tasks, a major challenge lies in handling datasets with imbalanced distributions, where one category is markedly more dominant than others. This disproportion might cause skewed educational outcomes, as standard classification approaches often lean towards the predominant group, leading to subpar identification of the less represented categories. Under such conditions, DRL stands out as a superior strategy for educating neural networks over conventional approaches. DRL addresses the problem of lopsided categorization by employing a system based on rewards [27]. Through carefully allocating incentives, it shifts the agent's attention towards instances belonging to the underrepresented categories, thus improving the detection of these rarer classes. The reward-centric model of DRL promotes a comprehensive decision-making process, prioritizing the discovery and classification of rare events or infrequently occurring categories.

In the realm of DRL, the primary goal of the agent is to select actions that optimize prospective benefits. The aggregation of rewards for forthcoming situations, symbolized by the reward value, gradually decreases over time, influenced by the discount rate $\gamma$, as illustrated in Eq. (2). In this formula, T corresponds to the concluding time-step of an episode [28].

$$R_t = \sum_{t'=t}^{T} \gamma^{t'-t} r_{t'} \tag{2}$$

where, $R_t$ represents the cumulative reward starting from time $t$, and $r_{t'}$ denotes the reward received at time $t'$. Q-values, representing the quality of state-action interactions, denote the anticipated outcome of policy $\pi$ upon executing action $a$ within state $s$. This is computed as depicted in Eq. (3).

$$Q^\pi(s,a) = E[R_t|s_t = s, a_t = a, \pi] \tag{3}$$

The most optimal action-value function, represented as the highest anticipated reward among all approaches after witnessing state $s$ and performing action $a$, is calculated as depicted in Eq. (4).

$$Q^*(s,a) = max_\pi E[R_t|s_t = s, a_t = a, \pi] \tag{4}$$

The function carries out the Bellman equation [29], which states that the supreme anticipated outcome for a particular maneuver is the sum of the benefits from the present maneuver and the utmost anticipated outcome from forthcoming maneuvers in the next instance. This concept is exemplified in Eq. (5).

$$Q^*(s,a) = E[r + \gamma \, max_{a'} Q^*(s',a')|s_t = s, a_t = a] \tag{5}$$

The computation of the ideal action-value function is methodically executed using the Bellman equation, as illustrated in Eq. (6).

$$Q_{i+1}(s,a) = E[r + \gamma \, max_{a'} Q_i(s',a')|s_t = s, a_t = a] \tag{6}$$

During the learning stage, as the network experiences state $s$, it generates a state-specific action. Subsequently, the system provides a reward r and transitions to the next state $s'$. These components are combined into a set $(s,a,r,s')$, subsequently stored in memory M. Groups of such sets, termed Batches B, are selected for performing gradient descent. The method for calculating loss is detailed in Eq. (7).

$$L_i(\theta_i) = \sum_{(s,a,r,s') \in B} (y - Q(s,a;\theta_i))^2 \tag{7}$$

Here, θ symbolizes the model's weights, while $y$ indicates the approximated objective for the Q function, evaluated as the summation of the reward linked with the state-action pair and the reduced maximum Q value in future instances, as illustrated in Eq. (8).

$$y = r + \gamma \, max_{a'} Q(s',a'; \theta_{k-1}) \tag{8}$$

It is important to recognize that the Q value assigned to the terminal state is initialized at zero. The gradient's magnitude for the loss function during the i-th iteration is ascertainable through Eq. (9).

$$\nabla_{\theta_i} L(\theta_i) = -2 \sum_{(s,a,r,s') \in B} (y - Q(s,a;\theta_i)) \nabla_{\theta_i} Q(s,a;\theta_i) \tag{9}$$

Through the execution of a gradient descent iteration on the loss function, adjustments are made to the model's weights under Eq. (10). This modification endeavors to lessen the discrepancy, where α denotes the learning rate dictating the extent of advancement within the optimization procedure.

$$\theta_{i+1} = \theta_i + \alpha \nabla_{\theta_i} Q(s,a;\theta_i) \tag{10}$$

*1) Problem formulation:* Within this paper, the application of the RL algorithm is directed towards the field of SHM. The ensuing explanation delineates the method's functioning and interpreting each component:

- State $s_t$: This matches the image captured at the temporal interval t. Here, an image refers to a graphical representation of the time series data. This image is composed of plots or graphs that visually depict the various anomalies identified in the acceleration time series of the bridge.

- Action $a_t$: The categorization executed on the image is regarded as an action. This signifies a choice carried out by the network, grounded in its prevailing comprehension of the objective.

- Action $r_t$: A reward is furnished for every categorization, designed to steer the network towards accurate categorization. The formulation of this remuneration process is expressed as:

$$r_t(s_t, a_t, y_t) = \begin{cases} +1 \,, a_t = y_t \text{ and } s_t \in D_O \\ -1 \,, a_t \neq y_t \text{ and } s_t \in D_O \\ \lambda \,, a_t = y_t \text{ and } s_t \in D_N \\ -\lambda \,, a_t \neq y_t \text{ and } s_t \in D_N \end{cases} \tag{11}$$

In this context, $D_N = \{Noraml\}$, and $D_O = \{Trend, Square, Missing, Minor, Drift, Outlier\}$ indicates the minority classes. Accurate or erroneous classification of a case from the prevalent category leads to an

incentive or penalty of +λ or -λ, respectively. This outlined method compels the network to focus on accurately identifying instances from the less frequent class by allocating a higher absolute value to the reward. Concurrently, the incorporation of the normal class and the flexible reward parameter within the range of 0<λ<1 adds complexity to the reward scheme. This allows for refined adjustment of the network's focus between the more and less prevalent classes.

## V. EMPIRICAL EVALUATION

In the meticulous assessment stage, a detailed and exhaustive analysis was carried out, contrasting our suggested model with six distinct deep-learning contenders, namely TransAnoNet [13], AnoSegNet [14], WaveletCNN [15], GAN-VAE [30], CNN-MIAD [31], and VibroCNN. This evaluation aimed to provide an all-encompassing insight into the strengths of our model vis-à-vis established methods. Moreover, we delved into different versions of our model by introducing three alternative variants for evaluation. The initial variant, termed "Proposed without dilated convolution," was based on a similar foundational architecture to our original model, yet did not incorporate dilated convolution. The subsequent variant, designated as "Proposed without RL," excluded the reinforcement learning component from the classification procedure. The third altered version, named "Proposed without DE," employed random initialization for the weights. We appraised these models using standard performance indicators, focusing specifically on measures like the F-measure and the geometric mean due to their proven effectiveness in tackling imbalanced datasets. The findings, detailed in Table II, resoundingly affirm the preeminence of our proposed model over all competing models, including those previously recognized as industry standards like AnoSegNet and TransAnoNet. Across every evaluative criterion, our model demonstrated consistent superiority over its rivals. Noteworthy accomplishments involve a marked error reduction exceeding 9% in the F-measure and surpassing 8% in the G-means indices. These notable advancements highlight the efficacy of our model in surmounting the difficulties associated with imbalanced datasets and its adeptness at furnishing more accurate forecasts. In juxtaposing our model with the modified iterations "Proposed without dilated convolution," "Proposed without RL," and "Proposed without DE," the indispensability of incorporating dilated convolution, RL, and DE methods becomes clear. Our model manifested a notable error rate reduction, approximately 5.35%, in comparison to its counterparts. This outcome accentuates the critical impact that the amalgamation of dilated convolution, RL, and DE strategies has in boosting the model's performance, thus solidifying their role as catalysts in the evolution of deep learning models.

In Fig. 2, we present the receiver operating characteristic (ROC) curves corresponding to the methodologies outlined in Table II. The area under the curve (AUC) serves as a pivotal metric for quantifying the performance of classifiers. An AUC score of 1 signifies impeccable discrimination ability, while a score of 0.5 suggests performance no better than random guessing.

It is worth highlighting that our proposed model emerges as the leader in this analysis, boasting a notable AUC value of 0.60. This solid outcome highlights its remarkable proficiency in accurately differentiating between favorable and unfavorable results, reinforcing the credibility of our approach as a potent predictive tool. Additionally, the "Proposed without RL" approach also demonstrates strong performance, achieving an AUC of 0.57, further affirming its ability to discern between positive and negative instances. In contrast, WaveletCNN and TransAnoNet, which achieve AUC scores of 0.46 and 0.49, respectively, offer less impressive performance. VibroCNN, GAN-VAE, and CNN-MIAD display even less favorable outcomes, with AUC values ranging from 0.43 to 0.45. Particularly, VibroCNN's meager AUC of 0.43, only slightly surpassing random chance, highlights its underwhelming performance. The ROC analysis vividly illustrates the varying degrees of performance among the evaluated methodologies. The exceptional predictive prowess demonstrated by our proposed method, whether in its standalone form or when coupled with RL, underscores the potency of our approach. Furthermore, it establishes a robust foundation for future enhancements and promising applications in the realm of predictive modeling, charting a path toward even more effective methodologies in the prediction domain. This remarkable performance positions our model as a key player in the field of predictive analytics.

TABLE II. EFFICIENCY INDICATORS OF THE SUGGESTED SYSTEM COMPARED TO RIVAL ADVANCED NETWORKS FOR SHM

|  | Accuracy | F-measure | G-means |
|---|---|---|---|
| TransAnoNet | 0.8104±0.0156 | 0.7802±0.1053 | 0.8102±0.0203 |
| AnoSegNet | 0.8001±0.2156 | 0.7371±0.0268 | 0.7902±0.1236 |
| WaveletCNN | 0.8005±0.2130 | 0.7202±0.0268 | 0.7906±0.2156 |
| GAN-VAE | 0.6801±0.0526 | 0.5402±0.1236 | 0.6405±0.0256 |
| CNN-MIAD | 0.7703±0.1563 | 0.6602±0.1036 | 0.7403±0.1265 |
| VibroCNN | 0.6704±0.0501 | 0.5501±0.0652 | 0.6462±0.0052 |
| Proposed without dilated convolution | 0.7904 ± 0.0517 | 0.7615 ± 0.1623 | 0.8014 ± 0.3622 |
| Proposed without RL | 0.8315 ± 0.0243 | 0.8200 ± 0.0417 | 0.8315 ± 0.0621 |
| Proposed without DE | 0.8425 ± 0.0123 | 0.8345 ± 0.0120 | 0.8436 ± 0.0505 |
| Proposed | 0.8601 ± 0.0384 | 0.8540 ± 0.0297 | 0.8760 ± 0.0123 |

Fig. 2. AUC chart for the suggested approach and alternative comparative techniques.

Fig. 3 showcases the confusion matrices for the proposed model, providing a detailed representation of its classification performance across different categories. From the matrix, we can observe the number of correct predictions (true positives) along the diagonal for each class, which are as follows: 13,271 for 'Normal', 5,363 for 'Trend', 2,608 for 'Square', 2,750 for 'Missing', 1,571 for 'Minor', 539 for 'Drift', and 409 for 'Outlier'. These figures suggest the model is most proficient at identifying the 'Normal' class and least proficient at identifying 'Outlier' instances, which could be because of their lower occurrence in the dataset. The off-diagonal numbers represent the instances where the model misclassified the inputs. For example, there are 76 instances where 'Normal' was incorrectly classified as 'Missing', and 80 instances where 'Square' was mislabeled as 'Normal'. Such misclassifications can diagnose and improve the model's performance, possibly by providing it with more representative training data or refining its feature detection capabilities.

Fig. 4 illustrates the evolution of error dynamics within the proposed model across 500 epochs. Commencing at an initial value of 12, the error undergoes a consistent descent as epochs unfold. This sustained decline signifies the model's progressive learning and enhancement of its predictive capabilities over time. It is significant to observe that the most pronounced decrease in error happens during the early training stages, slowly leveling off as the number of epochs increases. This trend indicates that with ongoing training, the rate of error reduction lessens, signifying a point where further error minimization from extended training becomes less impactful. Near the 425th epoch, a clear steadying of the error rate is observed, consistently hovering around a value of approximately 4.2962 in subsequent epochs. This leveling off of error rates suggests that continued training beyond this juncture is unlikely to result in notable enhancements in the model's forecasting accuracy. This stage may signal that the model has attained a state of convergence, reaching an accuracy level where additional fine-tuning might not bring considerable improvements. Alternatively, this stabilization could also suggest the emergence of overfitting concerns, especially if the model's performance on validation or test datasets ceases to improve. This insight into the error dynamics across epochs not only showcases the model's learning journey but also provides valuable guidance for fine-tuning training duration and preventing potential overfitting scenarios.

Fig. 3.    Comparative confusion matrices for the proposed model.



Fig. 4.    Comparative diagram of error dynamics.

## A.  Impact of the Reward Function

The allocation of rewards to both the more common and less frequent categories for accurate and inaccurate classifications is denoted by +1 and ±$\lambda$. The particular magnitude of $\lambda$ is determined by the ratio of frequent occurrences relative to rare events. As this ratio rises, it is expected that the ideal magnitude of $\lambda$ will diminish proportionally. In order to thoroughly investigate the influence of $\lambda$, we executed an extensive assessment of the suggested structure employing diverse $\lambda$ magnitudes, varying from 0 to 1 in steps of 0.1. Concurrently, the incentive for the more frequent category stayed unchanged. The detailed results are illustrated in Fig. 5. When adjusting $\lambda$ to 0, the effect of the dominant group turns negligible.

Conversely, with a value of $\lambda = 1$, both the more common and less common groups carry equivalent weight. The insights extracted from the analysis indicate that the framework achieves its peak effectiveness when $\lambda$ is established at 0.7, as observed across all assessed performance indicators. This observation suggests that the ideal $\lambda$ magnitude lies within the range of zero to one. It's important to recognize that although modulating $\lambda$ to reduce the impact of the dominant group is essential, configuring it too low might adversely affect the overall effectiveness of the entire structure. The evidence clearly indicates that the choice of $\lambda$ markedly affects the success of the structure. The suitable $\lambda$ magnitude depends on the comparative occurrences of more frequent and infrequent events, highlighting the need for careful determination to achieve the best results. This study underscores the intricate interplay between $\lambda$ and the framework's success, advocating for a balanced choice of $\lambda$ to strike a harmonious equilibrium between the two categories and foster effective results.

Fig. 5. Evaluation of the performance metrics of the proposed system under various settings of the parameter $\lambda$.

## B. Effect of Loss Function

The landscape of strategies available to combat the complexities arising from data imbalances in machine learning models is vast and diverse. It spans an array of techniques, ranging from the fine-tuning of data augmentation methods to the meticulous selection of aptly suited loss functions. The deliberate choice of an appropriate loss function plays a central role in ensuring the model's capacity to glean valuable insights from the underrepresented class embedded within the dataset. In our quest to unravel the nuances of the varying impacts of distinct loss functions, we embarked on a comprehensive exploration of five distinct contenders: WCE [32], BCE [33], DL [34], TL [35], and CL [36].

Among these contenders, both BCE and WCE have established themselves as widely adopted loss functions, treating positive and negative samples with equal significance. However, it's imperative to recognize that these functions might not be optimally configured to cater to datasets characterized by pronounced imbalances that accentuate the minority class. In stark contrast, DL and TL exhibit superior performance when confronted with skewed datasets, delivering

more favorable outcomes for the underrepresented class. Notably, CL emerges as a standout loss function, showcasing its prowess in scenarios where imbalanced data prevails. By skillfully adjusting the weights of the loss function, CL demonstrates its ability to prioritize intricate samples over simpler ones, thereby enhancing its adaptability in the face of challenging data distributions.

Our rigorous experimentation and analysis of these diverse loss functions are presented in meticulous detail in Table III. The outcomes unequivocally affirm the supremacy of CL over TL, leading to a substantial 3.72% reduction in the error rate concerning accuracy and an impressive 3.58% decrease in the F-measure. Nevertheless, it is crucial to underscore that, when benchmarked against the performance of our proposed model, CL exhibits a 1.5% deficit. These findings underscore the paramount significance of making a judicious selection of an appropriate loss function when navigating the intricacies of imbalanced data. Furthermore, they shine a spotlight on the commendable performance of our model in effectively addressing this prevalent and challenging issue in machine learning.

TABLE III. PERFORMANCE EVALUATION METRICS OF THE PROPOSED MODEL AGAINST VARIOUS LOSS FUNCTIONS IN SHM

|  | Accuracy | F-measure | G-means |
|---|---|---|---|
| WCE | 0.7303± 0.0269 | 0.7105± 0.1204 | 0.7412± 0.1120 |
| BCE | 0.7963± 0.0626 | 0.7536± 0.1203 | 0.7963± 0.0103 |
| DL | 0.7923± 0.0365 | 0.7821± 0.0056 | 0.8103± 0.1123 |
| TL | 0.8236± 0.2126 | 0.8023± 0.0145 | 0.8352± 0.1035 |
| CL | 0.8563± 0.0035 | 0.8325± 0.0032 | 0.8523± 0.0039 |

## C. Effect of CNNs

The architecture encompasses an array of CNNs that concurrently derive feature vectors from images. The quantity of CNNs utilized for feature extraction greatly influences the model's efficiency. Inadequate number of CNNs results in inadequate feature extraction, whereas too much may lead to problems such as overfitting or superfluousness. Both scenarios can diminish the model's overall utility. Therefore, carefully selecting the ideal number of CNN feature extractors is crucial. To identify this optimal number, we conducted a thorough and systematic analysis, evaluating the model's performance across a range of 1 to 7 CNN feature extractors. Our aim was to pinpoint the point where the model achieves peak functionality while maintaining a delicate balance between thorough feature extraction and operational efficiency. Our comprehensive experiments, as illustrated in Fig. 6, unequivocally demonstrate that three CNN feature extractors yield the model's best performance. Interestingly, as the number of CNNs increases, the model's performance declines, with six or seven extractors being less effective than using just one. This observed pattern highlights the existence of an optimal number of CNN feature extractors, maximizing the model's ability to capture relevant and discriminative features, ultimately leading to enhanced overall performance. Selecting three CNN feature extractors strikes a harmonious balance, optimizing the model's ability to extract essential information from input imagery, boosting its efficiency and effectiveness.



Fig. 6. Plotting the performance indicators of the suggested model while altering the quantity of convolutional feature extraction layers.

## D. Discussion

The proposed model signifies a significant advancement in the landscape of anomaly detection methods within the domain of SHM. By incorporating dilated convolutional and DE and RL techniques, this model showcases impressive predictive accuracy. These strides in technological innovation are especially pertinent considering the current challenges that the field of civil infrastructure encounters on a global scale.

However, it is essential to subject the model to critical examination within a broader context of its applicability. While the initial results are promising, they are inherently tied to data originating from a singular architectural marvel – a long-span cable-stayed bridge situated in China. While an in-depth focus on a specific dataset can yield valuable insights, it also presents the potential risk of confining the model to a narrow scope. Civil engineering marvels around the world encompass an immense range – from complex metro rail networks navigating urban mazes to towering skyscrapers reaching for the skies. Each of these structures is the culmination of distinct combinations of design, materials, and environmental factors, leading to unique challenges in structural health. For example, a dam nestled within mountainous terrain would encounter vastly different issues compared to a highway bridge spanning a saline estuary. Each structure reacts to external influences in a nuanced manner, whether it is the ceaseless battering of waves, vehicular loads, or the immense pressure of contained water. Therefore, while the anomalies identified in the Chinese bridge dataset offer invaluable insights, they might only scratch the surface of potential structural concerns when considering the full spectrum of potential issues.

Furthermore, alongside the diversity in structures, the environments in which they exist introduce an additional layer of complexity. The health of a structure isn't solely a reflection of its construction but also a result of its interactions with the environment [37]. From corrosion due to saline exposure to vibrations induced by seismic activities, the array of external stressors is extensive. This raises legitimate concerns about whether the proposed model, primarily trained on the Chinese bridge dataset, can seamlessly adapt to the myriad challenges that structures worldwide encounter. To address these concerns, several solutions can be implemented:

- Diverse Data Collection: Expanding the training dataset to include data from structures in different environmental conditions and geographic locations.

This would enhance the model's ability to generalize across a wide range of scenarios [38].

- Environmental Conditioning: Integrating environmental factors into the model, allowing it to learn how different environmental conditions affect structural health. This could involve adding parameters that account for local climate, pollution levels, and other relevant environmental data.

- Transfer Learning: Applying transfer learning techniques to adapt the model trained on the Chinese bridge dataset to other structures. This approach involves fine-tuning the model with smaller datasets from different structures, enabling it to adjust to new environments with minimal data.

- Regular Model Updates: Continuously updating the model with new data collected from various structures over time. This would ensure that the model stays relevant and effective in predicting structural health under changing environmental conditions [3].

- Hybrid Modeling Approaches: Combining the strengths of different modeling techniques, such as physics-based models and data-driven models [39]. This hybrid approach can leverage the accuracy of physics-based models in well-understood scenarios and the flexibility of data-driven models in complex, variable conditions.

- Real-time Environmental Monitoring: Integrating real-time environmental monitoring systems to provide continuous input to the model. This would allow the model to adjust its predictions based on current environmental conditions.

- Stress Testing and Simulations: Conducting stress tests and simulations under various environmental conditions to validate and improve the model's accuracy in different scenarios.

As the field of civil engineering advances, embracing new materials and groundbreaking construction techniques, the characteristics of potential structural irregularities are likely to transform. A cutting-edge SHM system must be proficient in identifying established problems and adept at signaling new, unexplored issues [40]. This capacity forms a crucial benchmark that the proposed model must meet. The implications are significant; failing to detect a key anomaly can result in catastrophic events, loss of human lives, and severe economic consequences. To enhance the proposed model's capability in this dynamic field, several approaches can be considered:

- Incorporation of Advanced Learning Algorithms: Utilizing machine learning and artificial intelligence algorithms that are capable of identifying patterns and anomalies not only from past data but also adapting to new trends. Techniques like unsupervised learning or deep learning can be particularly effective in recognizing unforeseen issues.

- Continuous Model Updating and Training: Regularly updating the model with the latest data from ongoing construction projects and newly developed materials. This will ensure that the model stays current and can recognize anomalies associated with new construction methodologies.

- Collaborative Data Sharing: Establishing a collaborative network with other civil engineering projects and research institutions for sharing data and insights. This collective approach can significantly broaden the spectrum of scenarios the model is exposed to, enhancing its ability to identify a wide range of anomalies.

- Predictive Analytics: Incorporating predictive analytics to forecast potential structural issues based on current trends and construction practices. This proactive approach can help in early identification and prevention of structural failures.

- Cross-Disciplinary Integration: Integrating knowledge from other fields such as materials science, meteorology, and environmental engineering. This interdisciplinary approach can provide a more comprehensive understanding of how various factors might contribute to new types of structural anomalies.

- Regular Sensitivity Analysis and Testing: Performing sensitivity analyses and stress tests under a variety of conditions to evaluate the model's effectiveness in detecting anomalies in different materials and construction methods.

- Expert Involvement and Feedback Loops: Engaging industry experts in regular reviews of the model's performance, ensuring that practical, real-world insights are incorporated. Establishing feedback loops can also aid in continuous improvement of the model.

As we chart our course ahead, multiple paths invite investigation. Firstly, testing the proposed model against a range of SHM datasets that include different types of structures, such as high-rise buildings, bridges, tunnels, and historical monuments, could reveal its extensive applicability [41]. Employing transfer learning techniques to adapt pre-existing models to these varied scenarios could be key in rapidly broadening the model's utility without necessitating extensive data gathering from each new structure type. In addition to these approaches, several other strategies could be beneficial:

- Cross-Functional Collaboration: Engaging with experts from different fields within civil engineering and data science to gain insights into specific structural characteristics and data processing techniques. This collaboration could enhance the model's accuracy and relevance across various structures.

- Real-Time Data Integration: Incorporating real-time monitoring data into the model to continually update and refine its predictive capabilities. This could include data from sensors monitoring weather conditions, material fatigue, and other relevant parameters.

- Customizable Model Parameters: Developing the model with customizable parameters that can be adjusted according to the specific requirements of different structures. This flexibility would allow for tailored applications, enhancing the model's effectiveness across diverse structural contexts.

- Scalability and Efficiency Improvements: Optimizing the model for scalability and computational efficiency to handle large datasets and enable its deployment in large-scale projects, such as city-wide infrastructure monitoring.

- Community Engagement and Feedback: Involving community feedback, especially from those who live or work in or near monitored structures, to provide ground-level insights into the model's performance and impact.

- Robust Validation and Testing: Conducting rigorous validation and testing under various conditions and scenarios to ensure the model's reliability and accuracy, particularly in critical and emergency situations.

- Policy and Regulatory Alignment: Ensuring that the model aligns with existing policies, standards, and regulatory requirements related to structural health and safety, to facilitate its acceptance and implementation.

Furthermore, the ever-changing characteristics of civil structures require that our SHM systems adapt and improve constantly. Implementing online learning paradigms in the proposed model would enable it to dynamically adjust to evolving structural health patterns. This could be achieved by continuously feeding the model with live data and allowing it to learn and update its parameters in real-time. The integration of diverse data sources, such as vibrations, strains, temperature changes, acoustic emissions, and even visual data from inspections, would significantly enrich the model's predictive accuracy [42]. Several additional steps can be taken to enhance the model's utility and efficiency:

- Edge Computing Implementation: Developing the model for deployment in edge computing environments where data processing occurs closer to the data source. This reduces latency and can be crucial for timely decision-making, especially in emergency scenarios.

- User-Friendly Interface Development: Creating intuitive user interfaces for the model that enable engineers and maintenance personnel to easily interpret and act upon the data and predictions provided by the system.

- Automated Alert and Reporting System: Integrating an automated system that generates alerts and detailed reports when anomalies are detected, thereby facilitating prompt and informed responses from the relevant authorities or maintenance teams.

- Interoperability with Existing Systems: Ensuring that the model is compatible with existing infrastructure management systems and can be seamlessly integrated into current workflows, enhancing its practicality and adoption.

- Regular Benchmarking and Validation: Regularly comparing the model's performance with other state-of-the-art anomaly detection systems in the field to validate its effectiveness and identify areas for improvement.

- Sustainability and Environmental Impact Assessment: Considering the environmental impact and sustainability of the model, especially in terms of its energy consumption and the materials required for sensor deployment and maintenance.

- Training and Education for Stakeholders: Providing comprehensive training and educational resources for engineers, technicians, and stakeholders to understand and effectively use the model in their operations.

Finally, it is important to delve deeper into the specific mechanisms and algorithms used in state-of-the-art techniques for handling imbalanced datasets in RL [43]. This involves examining different approaches, such as oversampling, under sampling, synthetic data generation, cost-sensitive learning, and novel reward shaping strategies [44, 45]. By contrasting these methods with our own, we can identify unique advantages or shortcomings in both theoretical and practical applications. Investigating how these techniques perform in diverse RL environments, ranging from simulated tasks to real-world applications, will provide a more holistic understanding of their adaptability and robustness. It would also be beneficial to explore the integration of our technique with other advanced machine learning strategies like deep learning, transfer learning, and meta-learning, to enhance its performance in handling imbalanced datasets. Such an in-depth analysis will not only fortify our research but also pave the way for future innovations in the field, fostering a more effective approach to tackling the challenges posed by imbalanced datasets in reinforcement learning [46].

## VI. CONCLUSION

This study introduced a groundbreaking model meticulously crafted to confront the intricate challenges associated with anomaly classification within SHM data. The proposed model harnessed a strategic fusion of dilated convolutional, RL, and DE techniques to achieve a high level of accuracy in its results. At its core, the model utilized a group of CNNs to extract essential feature vectors from input images concurrently. These extracted features were seamlessly integrated into downstream processes, bolstering the model's prowess in identifying complex patterns present in SHM data. The efficacy of the proposed model was rigorously validated through experimentation on an imbalanced dataset obtained from a long-span cable-stayed bridge in China-sourced from the IPC-SHM community. Handling imbalanced datasets poses distinct challenges in training classifiers, as the overrepresented class often exerts a disproportionate influence on the learning process, leading to suboptimal performance for the underrepresented class. To effectively address this concern, a novel approach was employed, integrating RL principles to formulate the training procedure as a series of interconnected

decisions. Within this framework, the dataset samples assumed the role of states, while the model operated as the agent, receiving appropriate rewards or penalties based on accurate or incorrect classifications, respectively. This adaptive strategy enabled the model to place a heightened focus on the underrepresented class, thereby enhancing classification outcomes. An innovative contribution to the training methodology was introduced by incorporating a mutation operator grounded in clustering principles within the framework of DE. This approach initiated the BP process by identifying a prominent cluster within the existing DE population. Subsequently, a novel update strategy was implemented to generate potential solutions, adding a layer of sophistication to the training process. The experimental results underscored the superior performance of the proposed model in the detection of multi pattern anomalies within SHM data, showcasing remarkable accuracy. Through the adept amalgamation of dilated convolutional, RL, and DE techniques, the model exhibited its potential as an advanced tool for anomaly detection within SHM systems. This capability is of utmost importance in safeguarding the structural integrity and safety of critical infrastructures, including vital components like bridges.

## REFERENCES

[1] J. P. Lynch, C. R. Farrar, and J. E. Michaels, "Structural health monitoring: technological advances to practical implementations [scanning the issue]," Proceedings of the IEEE, vol. 104, no. 8, pp. 1508-1512, 2016.

[2] W.-H. Hu, S. Said, R. G. Rohrmann, Á. Cunha, and J. Teng, "Continuous dynamic monitoring of a prestressed concrete bridge based on strain, inclination and crack measurements over a 14-year span," Structural Health Monitoring, vol. 17, no. 5, pp. 1073-1094, 2018.

[3] Z. Tang, Z. Chen, Y. Bao, and H. Li, "Convolutional neural network-based data anomaly detection method using multiple information for structural health monitoring," Structural Control and Health Monitoring, vol. 26, no. 1, p. e2296, 2019.

[4] H. Han, W.-Y. Wang, and B.-H. Mao, "Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning," in International conference on intelligent computing, 2005: Springer, pp. 878-887.

[5] I. Mani and I. Zhang, "kNN approach to unbalanced data distributions: a case study involving information extraction," in Proceedings of workshop on learning from imbalanced datasets, 2003, vol. 126: ICML, pp. 1-7.

[6] A. Fernández, S. García, M. Galar, R. C. Prati, B. Krawczyk, and F. Herrera, Learning from imbalanced data sets. Springer, 2018.

[7] S. Wang, W. Liu, J. Wu, L. Cao, Q. Meng, and P. J. Kennedy, "Training deep neural networks on imbalanced data sets," in 2016 international joint conference on neural networks (IJCNN), 2016: IEEE, pp. 4368-4374.

[8] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Learning deep representation for imbalanced classification," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 5375-5384.

[9] E. Lin, Q. Chen, and X. Qi, "Deep reinforcement learning for imbalanced classification," Applied Intelligence, vol. 50, pp. 2488-2502, 2020.

[10] S. V. Moravvej, S. J. Mousavirad, D. Oliva, and F. Mohammadi, "A Novel Plagiarism Detection Approach Combining BERT-based Word Embedding, Attention-based LSTMs and an Improved Differential Evolution Algorithm," arXiv preprint arXiv:2305.02374, 2023.

[11] T. Eltaeib and A. Mahmood, "Differential evolution: A survey and analysis," Applied Sciences, vol. 8, no. 10, p. 1945, 2018.

[12] S. Samudra, M. Barbosh, and A. Sadhu, "Machine Learning-Assisted Improved Anomaly Detection for Structural Health Monitoring," Sensors, vol. 23, no. 7, p. 3365, 2023.

[13] Q. Pan, Y. Bao, and H. Li, "Transfer learning-based data anomaly detection for structural health monitoring," Structural Health Monitoring, p. 14759217221142174, 2023.

[14] S. Li, L. Jin, Y. Qiu, M. Zhang, and J. Wang, "Signal anomaly detection of bridge SHM system based on two-stage deep convolutional neural networks," Structural Engineering International, vol. 33, no. 1, pp. 74-83, 2023.

[15] X. Ye, P. Wu, A. Liu, X. Zhan, Z. Wang, and Y. Zhao, "A Deep Learning-based Method for Automatic Abnormal Data Detection: Case Study for Bridge Structural Health Monitoring," International Journal of Structural Stability and Dynamics, p. 2350131, 2023.

[16] D. K. Green and A. Jaspan, "Applied Bayesian Structural Health Monitoring: inclinometer data anomaly detection and forecasting," arXiv preprint arXiv:2307.00305, 2023.

[17] R. Boccagna, M. Bottini, M. Petracca, A. Amelio, and G. Camata, "Unsupervised Deep Learning for Structural Health Monitoring," Big Data and Cognitive Computing, vol. 7, no. 2, p. 99, 2023.

[18] X. Lei, Y. Xia, A. Wang, X. Jian, H. Zhong, and L. Sun, "Mutual information based anomaly detection of monitoring data with attention mechanism and residual learning," Mechanical Systems and Signal Processing, vol. 182, p. 109607, 2023.

[19] Y. Yang and S. Nagarajaiah, "Data compression of structural seismic responses via principled independent component analysis," Journal of Structural Engineering, vol. 140, no. 7, p. 04014032, 2014.

[20] Y. Yang and S. Nagarajaiah, "Robust data transmission and recovery of images by compressed sensing for structural health diagnosis," Structural Control and Health Monitoring, vol. 24, no. 1, p. e1856, 2017.

[21] Y. Yang and S. Nagarajaiah, "Blind denoising of structural vibration responses with outliers via principal component pursuit," Structural Control and Health Monitoring, vol. 21, no. 6, pp. 962-978, 2014.

[22] S. Park, S. Kim, and J.-H. Choi, "Gear fault diagnosis using transmission error and ensemble empirical mode decomposition," Mechanical Systems and Signal Processing, vol. 108, pp. 262-275, 2018.

[23] Y. Bao, J. Li, T. Nagayama, Y. Xu, B. F. Spencer Jr, and H. Li, "The 1st international project competition for structural health monitoring (IPC-SHM, 2020): A summary and benchmark problem," Structural Health Monitoring, vol. 20, no. 4, pp. 2229-2239, 2021.

[24] J. Parra, L. Trujillo, and P. Melin, "Hybrid back-propagation training with evolutionary strategies," Soft Computing, vol. 18, no. 8, pp. 1603-1614, 2014.

[25] K. Deb, "A population-based algorithm-generator for real-parameter optimization," Soft Computing, vol. 9, pp. 236-253, 2005.

[26] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, and J. Pineau, "An introduction to deep reinforcement learning," Foundations and Trends® in Machine Learning, vol. 11, no. 3-4, pp. 219-354, 2018.

[27] M. Bahadori, M. Soltani, M. Soleimani, and M. Bahadori, "Statistical Modeling in Healthcare: Shaping the Future of Medical Research and Healthcare Delivery," in AI and IoT-Based Technologies for Precision Medicine: IGI Global, 2023, pp. 431-446.

[28] S. Danaei et al., "Myocarditis Diagnosis: A Method using Mutual Learning-Based ABC and Reinforcement Learning," in 2022 IEEE 22nd International Symposium on Computational Intelligence and Informatics and 8th IEEE International Conference on Recent Achievements in Mechatronics, Automation, Computer Science and Robotics (CINTI-MACRo), 2022: IEEE, pp. 000265-000270.

[29] E. Barron and H. Ishii, "The Bellman equation for minimizing the maximum cost," NONLINEAR ANAL. THEORY METHODS APPLIC., vol. 13, no. 9, pp. 1067-1090, 1989.

[30] J. Mao, H. Wang, and B. F. Spencer Jr, "Toward data anomaly detection for automated structural health monitoring: Exploiting generative adversarial nets and autoencoders," Structural Health Monitoring, vol. 20, no. 4, pp. 1609-1626, 2021.

[31] M. Zhao, A. Sadhu, and M. Capretz, "Multiclass anomaly detection in imbalanced structural health monitoring data using convolutional neural network," Journal of Infrastructure Preservation and Resilience, vol. 3, no. 1, p. 10, 2022.

[32] Ö. Özdemir and E. B. Sönmez, "Weighted cross-entropy for unbalanced data with application on covid x-ray images," in 2020 Innovations in

Intelligent Systems and Applications Conference (ASYU), 2020: IEEE, pp. 1-6.

[33] F. Huang, J. Li, and X. Zhu, "Balanced Symmetric Cross Entropy for Large Scale Imbalanced and Noisy Data," arXiv preprint arXiv:2007.01618, 2020.

[34] X. Li, X. Sun, Y. Meng, J. Liang, F. Wu, and J. Li, "Dice loss for data-imbalanced NLP tasks," arXiv preprint arXiv:1911.02855, 2019.

[35] S. S. M. Salehi, D. Erdogmus, and A. Gholipour, "Tversky loss function for image segmentation using 3D fully convolutional deep networks," in Machine Learning in Medical Imaging: 8th International Workshop, MLMI 2017, Held in Conjunction with MICCAI 2017, Quebec City, QC, Canada, September 10, 2017, Proceedings 8, 2017: Springer, pp. 379-387.

[36] S. A. Taghanaki et al., "Combo loss: Handling input and output imbalance in multi-organ segmentation," Computerized Medical Imaging and Graphics, vol. 75, pp. 24-33, 2019.

[37] A. Moallemi, A. Burrello, D. Brunelli, and L. Benini, "Model-based vs. data-driven approaches for anomaly detection in structural health monitoring: A case study," in 2021 IEEE International Instrumentation and Measurement Technology Conference (I2MTC), 2021: IEEE, pp. 1-6.

[38] Y.-M. Zhang, H. Wang, H.-P. Wan, J.-X. Mao, and Y.-C. Xu, "Anomaly detection of structural health monitoring data using the maximum likelihood estimation-based Bayesian dynamic linear model," Structural Health Monitoring, vol. 20, no. 6, pp. 2936-2952, 2021.

[39] C. Bigoni, "Numerical methods for structural anomaly detection using model order reduction and data-driven techniques," EPFL, 2020.

[40] Y. Bao, Z. Tang, H. Li, and Y. Zhang, "Computer vision and deep learning–based data anomaly detection method for structural health monitoring," Structural Health Monitoring, vol. 18, no. 2, pp. 401-421, 2019.

[41] Y. Zhang, Z. Tang, and R. Yang, "Data anomaly detection for structural health monitoring by multi-view representation based on local binary patterns," Measurement, vol. 202, p. 111804, 2022.

[42] X. Xu et al., "Anomaly detection for large span bridges during operational phase using structural health monitoring data," Smart Materials and Structures, vol. 29, no. 4, p. 045029, 2020.

[43] S. Susan and A. Kumar, "The balancing trick: Optimized sampling of imbalanced datasets—A brief survey of the recent State of the Art," Engineering Reports, vol. 3, no. 4, p. e12298, 2021.

[44] K. M. Hasib et al., "A survey of methods for managing the classification and solution of data imbalance problem," arXiv preprint arXiv:2012.11870, 2020.

[45] D. Ramyachitra and P. Manikandan, "Imbalanced dataset classification and solutions: a review," International Journal of Computing and Business Research (IJCBR), vol. 5, no. 4, pp. 1-29, 2014.

[46] J. M. Johnson and T. M. Khoshgoftaar, "Survey on deep learning with class imbalance," Journal of Big Data, vol. 6, no. 1, pp. 1-54, 2019.

# Application Effect of Human-Computer Interactive Gymnastic Sports Action Recognition System Based on PTP-CNN Algorithm

Yonge Ren[1], Keshuang Sun[2]*

Department of Physical Education, Jinzhong University, Jinzhong, 030600, China[1]
Department of Physical Education and Research, Fuzhou University, Fuzhou, 350108, China[2]

*Abstract*—With the rapid development of artificial intelligence technology, the recognition accuracy performance of traditional gymnastic sports action recognition system can no longer meet the needs of today's society. To address these problems, an improved action recognition algorithm combining Precision Time Protocal (PTP) and Convolutional Neural Networks (CNN) is proposed, and a human-computer interaction gymnastic action recognition system based on PTP-CNN algorithm is constructed. The performance test of the proposed PTP-CNN algorithm was conducted, and it was found that the accuracy of PTP-CNN algorithm was 92.8% and the recall rate was 95.2%, which was better than the comparison algorithm. The performance comparison experiments of the gymnastic action recognition system based on the PTP-CNN algorithm found that the recognition accuracy of the PTP-CNN gymnastic action recognition system was 96.3% and the running time was 3.4s, which was better than the other comparison systems. Comprehensive results can be found that the research proposed PTP-CNN recognition algorithm and improved gymnastic action recognition system can effectively improve the performance of traditional algorithms and models, which has practical application value and great application potential.

*Keywords—PTP; CNN; human-computer interaction; gymnastic sports; action recognition*

## I. INTRODUCTION

In recent years, with the development of technology, the application of artificial intelligence technology in the field of sports has become more and more extensive [1]. Among them, the human-computer interactive gymnastics sports action recognition system is a new type of artificial intelligence technology that can achieve automatic recognition and analysis of gymnastic actions through machine learning and deep learning algorithms [2]. The application of this system can improve the efficiency and quality of gymnastics training, which is of great significance to promote the development of physical education and training. At present, many scholars at home and abroad have researched and practiced the human-computer interactive gymnastics sports movement recognition system [3]. Some of these studies are based on machine learning algorithms and deep learning algorithms for design and implementation [4]. At present, the most commonly used gymnastics action recognition is mainly based on the convolutional neural network (Convolutional Neural Networks, CNN). The CNN can learn the feature representation automatically, and it has a strong non-linear

modeling capability. However, the traditional CNN mainly relies on the input of video frames and cannot make full use of the timing information of the action sequence. Moreover, the traditional CNN is prone to the problem of gradient disappearance or explosion when dealing with long time series. These problems limit the effectiveness of traditional sports systems in accuracy, real-time and personalized learning [5]. Accurate time protocol (PTP), as an accurate time synchronization protocol, has the advantages of high precision, high reliability and strong flexibility, which has a wide range in computer systems and networks. It can introduce time information in the action recognition task, and improve the accuracy and timing modeling ability of gymnastics movements of CNN by analyzing the forward and backward correlation in the action sequence [6]. At present, few studies have combined PTP with CNN and applied it in sports action recognition. Therefore, we propose to fuse PTP and CNN to build a PTP-CNN recognition algorithm and build a human-computer interactive gymnastic sports action recognition system based on it. It is expected to improve the efficiency and quality of gymnastics training and promote the development of physical education and training. Section I is the introduction to the article. The study describes the practical application of CNN and PTP in Section II. In Section III, PTP-CNN motion recognition algorithm and human-computer interactive gymnastics motion recognition system based on PTP-CNN are constructed. In Section IV, the action recognition algorithm and human-computer interaction action recognition system are tested. Results and discussion is given in Section V and Section VI concludes the paper.

## II. REVIEW OF THE LITERATURE

With the continuous in-depth research on CNN algorithms by domestic and foreign scholars, various CNN improvement models have been proposed and applied in several fields. In order to improve the recognition accuracy of coffee flowers, Wei et al. combined CNN model with binarization algorithm, selected a certain number of positive and negative samples from the original digital images for network model training, initially extracted coffee flowers based on the trained CNN model, and then further optimized its boundary information using binarization algorithm, and experimentally verified that the accuracy of this method for coffee flower classification was 93.7%, which has practical Application significance [7]. Chowdary et al. proposed a measurement system based on improved convolutional neural network in order to improve

the accuracy of mango leaf disease identification, and the proposed measurement was compared and analyzed with the system based on fuzzy algorithm, and the results showed that the accuracy of the proposed improved CNN-based system was 95.32, and this result indicated that the system could greatly improve the identification of mango leaf disease accuracy [8]. Joy and Vijayakumar found that the region-based convolutional neural network RCNN achieves good target detection accuracy but consumes more time on training and detection, so the proposed FAST RCNN algorithm with domain adaptive increments uses selective search to obtain the bounding box and feature extraction, thus overcoming the limitations of RCNN and improving the training and detection speed and accuracy [9]. Zhang's team proposed a deep learning model based on the combination of deep convolutional neural network and long and short-term memory network for the classification of arrhythmia intervals to address the problem that arrhythmias are difficult to diagnose accurately. Ten-fold cross-validation of the method showed that the average accuracy of the method was 99.06%, which is of great significance in clinical applications [10]. Zhao et al. For the problem of distortion and artifacts in lossy compressed video, a learning model with variable filter size residuals is proposed based on CNN algorithm, and the effectiveness of this model is measured using a combination based on color sensitivity, and extensive experimental results show that it has a better performance than existing methods in terms of efficiency improvement after video coding. [11].

With the rapid development of information technology, there are more and more methods applied in the field of action recognition. Rubin's team proposed a faster region based convolutional neural network structure to address the problem of difficult to accurately recognize real-time gestures, and tested its performance on a standard data set. Carvalho et al. proposed a multi-standard action-based human-robot interaction framework in order to improve the accuracy of socially assisted robot action recognition. [12]. Carvalho et al. tested the method offline and online, and the results showed that the accuracy of the method exceeded 96.7%, which is practical and can be used in educational proposals. [13]. Hu's team addressed the problem that current action recognition methods tend to ignore the reversibility of skeleton data in the temporal dimension [14]. Gao et al. proposed a new forward-inverse adaptive graph convolutional network for skeleton-based action recognition to address the problem that the current graph convolutional network models focus more on spatial information and ignore temporal information, and empirically analyzed the method. Gao et al. propose a unified attention model that integrates channel, space, and time, and the model is tested for performance and found to have the best performance compared to other similar action recognition methods [15].

The above studies fully illustrate that the CNN improvement model has been widely used in several fields, and there are also various methods applied in the field of action recognition. However, there are fewer studies combining PTP methods with CNN algorithms, so the study combines PTP methods with CNN algorithms to obtain PTP-CNN algorithms, and applies the improved algorithms to

human-computer interactive gymnastic sports action recognition, expecting to improve the accuracy of gymnastic sports action recognition in this way and promote the further development of gymnastics course intelligence.

## III. CONSTRUCTION OF HUMAN-COMPUTER INTERACTION GYMNASTIC SPORTS ACTION RECOGNITION SYSTEM BASED ON PTP-CNN ALGORITHM

### A. CNN Action Recognition Algorithm Based on the Principle of PTP Protocol

PTP protocol, also known as IEEE1588 protocol, is currently a mainstream time synchronization system, which is perfectly suitable for modern communication technology at the same time, but also well combined with computer hardware time equipment [16]. PTP can be placed in the computer network multiple clocks, and the master clock source and other clock sources in the form of telegrams time, to achieve accurate synchronization of network time. The PTP protocol clock type is divided into ordinary clock, boundary clock, end-to-end transparent clock and point-to-point transparent clock [17]. PTP synchronization principle implementation process is shown in Fig. 1.



Fig. 1. Implementation process of PTP synchronization principle.

As it can be seen from Fig. 1, the PTP message flow is divided into four modules. First, the host sends an ANNOUNCE message to all devices in the slave, and the devices in the slave listening state will receive the ANNOUNCE message and set the host's clock source to the best master clock and set it to the uncalibrated state. Subsequently, the host sends a SYNC message to the slave and records the timestamp $t_1$, the slave receives it and records the timestamp $t_2$. The third module is the host sends a FOLLOW-UP message to the slave with the timestamp $t_1$ and the slave receives it and sends a DELAY-REQ message back to the host with the timestamp $t_3$. The host receives the

DELAY-REQ message and records the timestamp $t_4$. Finally, the host sends the DELAY-RESP message and the timestamp $t_4$ to the slave, and the slave finally gets a total of four timestamps. The slave obtains the link delay between the host and the slave by calculating the time deviation of the obtained timestamps from that of the host. The four timestamps $t_1$, $t_2$, $t_1$ and $t_4$ obtained by the slave are calculated as shown in Eq. (1).

$$\begin{cases} t_2 = t_1 + delay_{MS} + offset \\ t_4 = t_3 + delay_{SM} - offset \end{cases} \tag{1}$$

In Eq. (1), $offset$ is the time deviation from slave to host; $delay_{MS}$ is the link delay from host to slave; $delay_{SM}$ is the link delay from slave to host. When the values of $delay_{MS}$ and $delay_{SM}$ are the same, the link delay and time deviation are calculated as shown in Eq. (2).

$$\begin{cases} delay_{SM} = \dfrac{(t_2 - t_1) + (t_4 - t_3)}{2} \\ offset = t_2 - delay_{SM} \end{cases} \tag{2}$$

In Eq. (2), $t_1$, $t_2$, $t_1$ and $t_4$ are the time stamps obtained by the slave; $offset$ is the time deviation between the slave and the host; $delay_{MS}$ is the link delay from the host to the slave; $delay_{SM}$ is the link delay from the slave to the host. Due to the problem of link delay jitter in time synchronized networks in real networks. The study uses the second-order Kalman filtering algorithm to process the time deviation and link delay. The formula of first-order exponential smoothing filtering at this point is shown in Eq. (3).

$$y(n) = \alpha \cdot x(n) + (1-\alpha) \cdot y(n-1) \tag{3}$$

In Eq. (3), $y(n)$ denotes the value after the first $n$ filtering and the initial value is 0; $\alpha$ is the coefficient of 0-1; $x(n)$ denotes the observed value of the first $n$. The cutoff frequency of the exponential smoothing filter is calculated as shown in Eq. (4).

$$\begin{cases} Y(Z) = \alpha \cdot X(Z) + (1-\alpha) \cdot Y(Z) \cdot z^{-1} \\ H(Z) = \dfrac{Y(Z)}{X(Z)} = \dfrac{\alpha}{1-(1-\alpha) \cdot z^{-1}} \end{cases} \tag{4}$$

In Eq. (4), $H(Z)$ denotes the system function; $Y(Z)$ denotes the discrete $Z$ transformation of $y(n)$; $\alpha$ is the coefficient of 0-1; $Y(Z)$ denotes the discrete $Z$ transformation of $x(n)$. The system frequency response is calculated as shown in Eq. (5).

$$H(e^{j\omega}) = H(Z)\big|_{z=e^{j\omega}} = \dfrac{\alpha}{1-1(1-\alpha) \cdot e^{-j\omega}} \tag{5}$$

In Eq. (5), $H(Z)$ denotes the system function; $\alpha$ is the coefficient of 0-1. The system amplitude and frequency response is shown in Eq. (6).

$$\left| H(e^{j\omega}) \right| = \left| \dfrac{\alpha(\cos\omega + j\cdot\sin\omega)}{\cos\omega - (1-\alpha) + j\cdot\sin\omega} \right| = \dfrac{\alpha}{\sqrt{1+(1-\alpha)^2 + j\cdot\sin\omega}} \tag{6}$$

In Eq. (6), $\alpha$ is a factor of 0-1. $_{-3dB}$ When the system cutoff frequency is calculated, the formula is shown in Eq. (7)

$$f = \dfrac{\arccos\left(1 - \dfrac{\alpha^2}{2(1-\alpha)}\right)}{2\pi} \tag{7}$$

In Eq. (7), $\alpha$ is a factor of 0-1. Since exponential smoothing filtering can effectively reduce the time deviation and thus improve the degree of system stability. The equation of control filtering at this point is shown in Eq. (8)

$$\begin{cases} l(n) = k(n) \cdot K_i + l(n-1) \\ m(n) = k(n) \cdot K_p + l(n) \end{cases} \tag{8}$$

In Eq. (8), $l(n)$ is the output value of the $n$ integral control; $k(n)$ is the observed value of $n$; $k_i$ is the coefficient of the integral control; $k(n)$ is the output value of the $n$ exponential filter control; and $k_p$ is the coefficient of the proportional control. The filtering process is shown in Eq. (9).

$$\begin{cases} filter_{es} = \alpha \cdot offset(n) + (1-\alpha) \cdot filter_{es}(n-1) \\ filter_l(n) = filter_{es}(n) \cdot K_l + filter_l(n-1) \\ filter_p(n) = filter_{es}(n) \cdot K_p + filter_p(n-1) \end{cases} \tag{9}$$

In Eq. (9), $filter_{es}$ is the exponential smoothing filter output value of $n$; $\alpha$ is the smoothing filter coefficient between 0 and 1; $offset(n)$ is the time deviation before the filter of $n$; $filter_i(n)$ is the integral control output value of $n$ of the filter function; $K_l$ is the integral control coefficient; $filter_p(n)$ is the proportional control output value of $n$ of the filter function; $K_p$ is the proportional control coefficient. After implementing the Kalman filter-based clock taming, the CNN action recognition algorithm based on the PTP principle needs to be constructed where the basic structure of CNN as shown in Fig. 2.

Fig. 2. CNN basic structure.

As shown in Fig. 2, CNN has five layers structure, which are input layer, convolutional layer, pooling layer, fully connected layer and output layer. Among them, the CNN convolutional layer is used for recognition and feature extraction of input data, including multiplication and addition of matrices, and its operation is closely related to the convolutional kernel [18-19]. The convolutional kernel is a feature extractor with sparse connectivity and weight sharing in the convolutional layer, and the convolutional kernel is trained to output a feature set that meets the needs [20]. The operation formula of the convolution layer is shown in Eq. (10).

$$x^l_j = f\left(\sum^n_{i \in M_j} x^{l-1}_i * k^l_{i,j} + b^i_j\right) \tag{10}$$

In Eq. (10), $x^l$ and $f$ are the output and activation function of the $l$ layer, respectively; $x^{l-1}_i$ is the output of the $l-1$ layer; $k^l$ and $b^i_j$ are the convolution kernel and offset term of the $l$ layer; and $M_j$ is the selected input feature set. The pooling layer divides the obtained feature set and reduces the dimensionality of the features to reduce the computational effort and enhance the robustness. The common operations of the pooling layer are mainly maximum pooling and mean pooling, and the two pooling methods are shown in Eq. (11).

$$a = \begin{bmatrix} 2 & 3 & 0 & 3 \\ 1 & 4 & 4 & 3 \\ 5 & 6 & 4 & 3 \\ 1 & 0 & 0 & 1 \end{bmatrix}, a_{max} = \begin{bmatrix} 4 & 4 \\ 6 & 4 \end{bmatrix}, a_{ave} = \begin{bmatrix} 2.5 & 2.25 \\ 3 & 2 \end{bmatrix} \tag{11}$$

Finally, in order to prevent overfitting or underfitting of the model, reasonable optimization of the parameters is required. The study can use cross-entropy loss function and back propagation algorithm to calculate the error in fault multi-classification diagnosis. The formula of cross-entropy loss function is shown in Eq. (12).

$$H(p,q) = -\sum_x p(x)\log q(x) \tag{12}$$

In Eq. (12), $p(x)$ is the target distribution; $q(x)$ is the predictive distribution. The backpropagation algorithm is an algorithm for training a feedforward neural network for a given input pattern with a known classification using the chain derivation method. The backpropagation algorithm is the most

common and effective method for training artificial neural network algorithms, and the essence is the error between the output and the target, as shown in Eq. (13).

$$\delta^{l-1} = \frac{\partial J}{\partial z^{l-1}} = \frac{\partial J}{\partial z^l}\frac{\partial z^l}{\partial z^{l-1}} = \delta^l \frac{\partial z^l}{\partial z^{l-1}}\frac{\partial a^{l-1}}{\partial z^{l-1}} \tag{13}$$

In Eq. (13), $\delta^l$ is the error of the objective function $J$ to $z^l$. Finally, the PTP principle is fused with the CNN algorithm to construct the PTP-CNN recognition algorithm. the structure of the PTP-CNN recognition algorithm is shown in Fig. 3.



Fig. 3. PTP-CNN recognition algorithm structure.

Fig. 3 shows the structure of the PTP-CNN recognition algorithm, and the solid arrows in the figure represent the back-and-forth relationship between the modules in the PTP-CNN algorithm. As shown in Fig. 3, the PTP-CNN algorithm can extract the spatial features of video actions through CNN networks and fuse the convolutional features of different levels of CNN networks to enhance the feature representation. The temporal information in the video frames is modeled by PTP, and finally the classification results are obtained in using the class ware to recognize the human actions in gymnastics videos. To improve the recognition efficiency, the study first preprocesses the gymnastic video data with video frames, and adjusts the video frame size to (224*224) as the input. After completing the video pre-processing, the spatial features in the action video are extracted by CNN network and segmented into (224*224*3) size as input, the "3"in (224*224*3) indicates the channel

dimension size. After that, the PTP-CNN algorithm fuses the shallow features with the deep features in the video through pooling and splicing operations to compensate for the missing feature information such as location contours. After several convolutions and pooling, the feature map with the output size of (7*7*2048) is pooled globally on average and expanded to (1*2048), using dropout to avoid overfitting of the results. Finally, the extracted feature vectors are input to the PTP time synchronization server to model the temporal information of the actions and the classification of the actions is achieved by a Softmax classifier.

### B. Construction of Gymnastic Movement Recognition System Based on PTP-CNN Algorithm

After completing the construction of PTP-CNN recognition algorithm, the research will develop a gymnastic action recognition system. Since visual information is easily disturbed by external factors in real scenes, the accuracy and robustness of action recognition relying only on a single operational logic is very poor, so the research proposes a gymnastic action recognition system based on PTP-CNN algorithm. At the same time, in order to avoid the problem that the structure of the gymnastic action recognition system is confusing and difficult to expand, the research firstly designs the structure of the gymnastic action recognition system. The basic structure of the gymnastic action recognition system based on PTP-CNN algorithm is shown in Fig. 4.

As shown in Fig. 4, the gymnastic action recognition system proposed in the study takes reliability, practicality and scalability as design principles, and divides the system structure into four modules: web application layer, business logic layer, software development layer and basic function layer. The web application layer is the display page of the system to users, which includes user registration and login, video upload, action classification and recognition query functions. The business logic layer is the functional basis of the application layer, the main task is to manage the user database, complete the video pre-processing operation, and the video recognition and classification result analysis by PTP-CNN algorithm. The software development layer is the tool layer for system construction, including Pytorch framework, Python language, HTML/CSS front-end page development language and lightweight Flask framework, etc. The basic platform layer is the platform for the operation of the PTP-CNN gymnastic movement recognition system, and all the functions in the system cannot be realized without the support of hardware equipment and operation system. The hardware environment configuration of the gymnastic movement recognition system is AMD R5-4600H, 3.0GHz six-core twelve-thread processor, NVIDIA GeForce GTX1650Ti graphics card, 16GB running memory, 512GB solid state drive. The software environment is configured with Python language and Pytorch deep learning framework. the principle of PTP-CNN gymnastic action recognition system is shown in Fig. 5.



Fig. 4.   Basic structure of gymnastic movement recognition system based on PTP-CNN algorithm.



Fig. 5.   Principle of PTP-CNN gymnastic movement recognition system.

Fig. 6.    The overall framework of PTP-CNN gymnastics motion recognition system.

As shown in Fig. 5, the proposed PTP-CNN gymnastic action recognition system consists of three modules, namely, the data pre-processing module, the human-computer interaction module and the PTP protocol engine module. In the data pre-processing module, the system decomposes the video in the dataset into continuous frames and feeds them into the HCI module. The principle of HCI module is to finish the feature recognition and feature extraction of gymnastics video by CNN algorithm. Finally, the PTP protocol engine module takes the video action information and different action features extracted by the CNN algorithm, models and fuses them in 3D, and outputs the gymnastic action feature recognition results. The overall framework of PTP-CNN gymnastics action recognition system is shown in Fig. 6.

As shown in Fig. 6, the PTP-CNN gymnastic action recognition system is mainly divided into the user's login and registration module, and the action recognition module of the underlying logic of the system. In the login and registration module, after the user enters the system login window, the system takes the user's username at the time of registration as the login account, and judges the submitted username when the user submits the registration information. If the user name is recognized to exist in the database, the user registration fails, otherwise the registration is successful, at this time the PTP-CNN gymnastics action recognition system generates a unique user number primary key for the subsequent query operation of the user. After the registration is completed, users can enter the system by entering the correct account password, at which time they can upload the gymnastics video. Since the PTP-CNN gymnastics action recognition system only supports video uploads in avi and mp4 formats, the system will identify the format of the video uploaded by the user. If the format is correct, it will enter the video preview stage; if the video format is wrong, it will return to the video upload stage. After entering the video recognition stage, the system will call the PTP-CNN algorithm in the background for action recognition and write the recognition result to the database history. After finishing the gymnastic action classification recognition, the background will package the recognition results into JSON format data to return to the front-end, and the results will be displayed in the user interface.

## IV. EMPIRICAL ANALYSIS OF HUMAN-COMPUTER INTERACTIVE GYMNASTIC SPORTS ACTION RECOGNITION SYSTEM BASED ON PTP-CNN ALGORITHM

### A. Analysis of the Effectiveness of PTP-CNN Recognition Algorithm

To verify the effectiveness of the PTP-CNN recognition algorithm for gymnastic sports action recognition, the study uses the public dataset UTKinect-Action3D to validate the effectiveness of the PTP-CNN recognition algorithm. the UTKinect-Action3D dataset contains activities corresponding to the gymnastic action recognition system, such as finishing exercises, stretching exercises, chest expansion exercises, full body movement, body rotation movement and jumping movement. The PTP-CNN recognition algorithms are compared and analyzed by the recognition accuracy, recall, F1 value, verification loss value (val-loss), and verification accurate values (val-acc) of these six corresponding actions. (Visual Geometry Group Network-16, VGG16), Residual Network (ResNET) and Dynamic Time Warping (DTW) algorithms. The recognition accuracy and recall curves of the compared algorithms are shown in Fig. 7.

Fig. 7(a) shows the recognition accuracy curves of the compared algorithms. From Fig. 7(a), it can be seen that the recognition accuracy of each algorithm increases with the number of iterations, and the PTP-CNN algorithm proposed in the study has an overall higher accuracy than the other algorithms, with an accuracy rate of up to 94.3% and an average accuracy rate of 92.8%. Fig. 7(b) shows the recall curves of the compared algorithms. From Fig. 7(b), it can be seen that the recall rate of each algorithm is smooth and does not change with the number of experiments, among which the PTP-CNN algorithm proposed in the study has a higher recall rate than the other algorithms, and its average recall rate is 95.2%. From the above results, it is clear that the accuracy performance and recall performance of the PTP-CNN algorithm proposed in the study are better than the other algorithms. Fig. (8) shows the accuracy-recall rate curves and F1 value comparison results of each compared algorithm.

(a)Accuracy curve of comparison algorithm                    (b)Comparison algorithm recall curve

Fig. 7.   Accuracy and recall curves of each algorithm.



(a)The first comparative experiment                    (b)Compare the F1 score of the algorithm

Fig. 8.   Accuracy-recall and F1 score curves of each algorithms.



(a)Val-loss curves of various comparison algorithms                    (b)Val-acc curves of various comparison algorithms

Fig. 9.   Validation losses and validation accuracy of each algorithm.

Fig. 8(a) shows the accuracy-recall curves of the compared algorithms. From Fig. 8(a), it can be seen that the accuracy-recall curve of the proposed PTP-CNN algorithm has the largest area under the line of 0.81 compared with other algorithms, which is 0.5 larger than the accuracy-recall curve of CNN. Fig. 8(b) shows the F1 values of each comparison algorithm. From Fig. 8(b), it can be seen that the proposed PTP-CNN algorithm has the largest F1 value of 0.93 compared to the other compared algorithms, which is 0.08 higher than the F1 value of CNN. In summary, it can be seen that the proposed PTP-CNN algorithm has the best performance in terms of accuracy-recall rate and F1 value performance. Fig. 9 shows the val-loss and val-acc values of

the compared algorithms.

Fig. 9(a) shows the val-loss curves of each comparison algorithm. From Fig. 9(a), it can be seen that the val-loss curve of the proposed PTP-CNN algorithm has the lowest overall val-loss curve and the smallest fluctuation compared to the other comparison algorithms, with an average val-loss value of 0.72 and a fluctuation of 0.31. Fig. 9(b) shows the val-acc values of each comparison algorithm. From Fig. 8(b), it can be seen that the PTP-CNN algorithm proposed in the study has the largest val-acc value of up to 1.51 compared with the other comparison algorithms, which is 0.58 higher than the highest val-acc value of CNN. In summary, the results show that the PTP-CNN algorithm proposed in the study has

the best performance in terms of the performance of verification loss and verification accuracy.

*A. Comparison Experiment of Gymnastic Movement Recognition System Based on PTP-CNN Algorithm*

To test the recognition performance of the PTP-CNN algorithm-based gymnastic movement recognition system, the study conducts a comparative performance analysis of the system. The study tests the performance of the system by comparing the recognition accuracy, recognition error and system running time, etc. The comparison system is a gymnastic action recognition system based on ResNET, VGG16 and CNN algorithm. The system test platform is composed of KinectV1.0, Windows10, VisualStudio and Unity. The recognition accuracy of each system is shown in Fig. 10.



Fig. 10. Identification accuracy of each system.

Fig. 10 shows the recognition accuracy of the action recognition system based on PTP-CNN, ResNET, VGG16 and CNN algorithms for stretching, chest expansion, body turn, jumping, full-body and finishing movements. As shown in Fig. 10, the recognition accuracy of the PTP-CNN algorithm-based gymnastics action recognition system is higher than other comparison systems in six movements as a whole, among which the system has the lowest recognition accuracy of 93.4% for finishing movements and the highest recognition accuracy of 98.5% for full-body movements. The average recognition accuracy of the PTP-CNN algorithm-based gymnastic movement recognition system is 96.3%, which is 10.1% more accurate than that of the ResNET algorithm-based recognition system. Summing up the results, it can be concluded that the gymnastic movement recognition system based on PTP-CNN algorithm proposed in the study has the best performance in terms of movement recognition accuracy. The confusion matrix of the PTP-CNN algorithm-based gymnastic action recognition system is shown in Fig. 11 when the recognition results of the actions are compared with the actual actions.

Fig. 11 shows the confusion matrix of the gymnastic action recognition system based on PTP-CNN algorithm. From Fig. 11, it can be seen that the recognition accuracy of the gymnastic action recognition system proposed in the study is high on all six actions, indicating that the system has excellent recognition and classification ability on and distinguished actions, which has practical use value. To further verify the practical use performance of the PTP-CNN algorithm-based gymnastic action recognition system, the study conducted empirical experiments on the system to analyze its system recognition error and system computing speed, and the results of the empirical experiments are shown in Fig. 12.



Fig. 11. PTP-CNN gymnastics recognition system confusion matrix.

Fig. 12. Identification error and calculation time of each system.

Fig. 12(a) shows the recognition errors of gymnastic movements for each comparison system. From Fig. 12(a), it can be seen that among the four comparison systems, the system with the largest recognition error is the CNN algorithm-based gymnastic movement recognition system, whose recognition error is 11.5%. And the recognition system based on PTP-CNN algorithm has the smallest recognition error among the comparison systems, which is 2.4% and 9.1% lower than the gymnastic action recognition system based on CNN algorithm. Fig. 12(b) shows the operation speed of each comparison system. From Fig. 12(b), it can be seen that among the four comparison systems, the longest running speed is the CNN algorithm-based gymnastic movement recognition system, which runs for 10.9 s. The PTP-CNN algorithm-based recognition system has the shortest running time among the comparison systems, which is 3.4 s and 7.5 s shorter than the CNN algorithm-based gymnastic movement recognition system. Gymnastic movement recognition system has the best performance in terms of recognition error and running speed, and has significantly improved the performance compared with the traditional gymnastic movement recognition system.

## V. RESULTS AND DISCUSSION

It is found that the PTP-CNN motion recognition system can accurately identify the stretching, chest expansion, body turning, body jumping, whole body movement and finishing movement in gymnastics, with the accuracy rate of 96.3%, which is better than the traditional CNN motion recognition system. The findings are consistent with Yu et al. [21], which improved the accuracy of action classification and were applied to EMG control. In addition, the study also found that PTP-CNN human-computer interactive gymnastics motion recognition system has a short running time and has practical application value. Similar to the study of Majd et al. [22]., can be applied in the fields such as video surveillance. In addition, the action recognition system was applied to the path planning of the traffic system by Chen et al. [23]. Therefore, the proposed human-computer interaction gymnastics motion recognition system has good application potential in the fields of medicine, video surveillance, traffic and sports. In the medical field, the system can be used for rehabilitation training and evaluation, helping doctors to monitor patients'

exercise recovery and provide targeted rehabilitation programs. In the field of video surveillance, the system can be used to monitor and identify human movements in real time, and help security personnel quickly detect abnormal or criminal behaviors. In the field of traffic, the system can be used to identify the driver's action and posture, monitor the driver's fatigue driving situation, remind the driver to pay attention to safety, so as to reduce the occurrence of traffic accidents. Most importantly, in the field of sports, the system can be applied in training and competition to help coaches and athletes analyze and improve movement skills, and improve training results and competition performance. In addition, the system can be used to evaluate the performance of the players and provide objective basis for the judges. In conclusion, the human-computer interactive gymnastics motion recognition system has wide application potential to play an important role in medicine, video surveillance, transportation and sports.

## VI. CONCLUSION

In order to improve the recognition accuracy and operation speed of the traditional gymnastic action recognition system, the study proposes an action recognition algorithm that integrates the PTP principle and CNN algorithm, and builds a gymnastic action recognition system based on the PTP-CNN algorithm based on it. The performance tests of the proposed PTP-CNN algorithm and the improved gymnastic action recognition system are conducted. The results show that the PTP-CNN algorithm has the highest accuracy of 94.3%, the average accuracy of 92.8%, and the recall rate of 95.2%, which are better than the rest of the comparison algorithms in terms of accuracy performance and recall rate performance. In addition, the study also conducted performance comparison experiments on the improved gymnastic movement recognition system based on the PTP-CNN algorithm. The results show that the recognition accuracy of PTP-CNN gymnastic action recognition system is 96.3%, which is better than the other comparison systems. In addition, it is found that the average running time of PTP-CNN gymnastic action recognition system is 3.4s, which is lower than other comparison systems. The above results can be found that the proposed PTP-CNN recognition algorithm and gymnastic movement recognition system are better than the comparison algorithm and system in terms of recognition accuracy and

running speed, and have practical application value. However, the study also has some limitations, compared to the gymnastics posture, including the extension, distortion and rotation of the body. This attitude change poses certain challenges to the accuracy and robustness of the algorithm. Future studies can explore more effective posture feature extraction methods according to the large posture changes in gymnastics movements, and improve the algorithm's ability to identify and model posture changes.

## REFERENCES

[1] Qian J. Research on Artificial Intelligence Technology of Virtual Reality Teaching Method in Digital Media Art Creation. Journal of Internet Technology, 2022, 23(1):125-132.

[2] Modi N, Singh J. Role of Eye Tracking in Human Computer Interaction. ECS transactions, 2022, 107(1):8211-8218.

[3] Xu W. From Automation to Autonomy and Autonomous Vehicles: Challenges and Opportunities for Human-Computer Interaction. interactions, 2021, 28(1):48-53.

[4] Oslund S, Washington C, So A, Chen T, Ji H. Multiview Robust Adversarial Stickers for Arbitrary Objects in the Physical World. Journal of Computational and Cognitive Engineering, 2022, 1(4): 152-158.

[5] Bae J, Lee D H. PTP Tracking Scheme for Indoor Surveillance Vehicle by Dual BLACM With Hall Sensor. IEEE Transactions on Industry Applications, 2022, 58(4):5238-5247.

[6] Nimrah S, Saifullah S. Context-Free Word Importance Scores for Attacking Neural Networks. Journal of Computational and Cognitive Engineering, 2022, 1(4): 187-192.

[7] Wei P, Jiang T, Peng H, Jin H, Huang H. Coffee Flower Identim fication Using Binarization Algorithm Based on Convolutional Neural Network for Digital Images. Plant Phenomics, 2020, 5:101-115.

[8] Chowdary M S, Puviarasi R. Accuracy Improvement in Disease Identification of Mango Leaf using CNN Algorithm Compared with Fuzzy Algorithm. ECS transactions, 2022, 107(1):11889-11903.

[9] Joy F, Vijayakumar V. Multiple Object Detection In Surveillance Video With Domain Adaptive Incremental Fast Rcnn Algorithm. Indian Journal of Computer Science and Engineering, 2021, 12(4):1018-1026.

[10] Zhang P, Hang Y, Ye X, Guan P, Hu W. A United CNN-LSTM Algorithm Combining RR Wave Signals to Detect Arrhythmia in the 5G-Enabled Medical Internet of Things. IEEE Internet of Things Journal, 2021, 9(16):14563-14571.

[11] Zhao H, He M, Teng G, Shang X, Wang G, Feng Y. A CNN-Based Post-Processing Algorithm for Video Coding Efficiency Improvement. IEEE Access, 2020, 8(1):920-929.

[12] Rubin B S, Kumar V S. An Efficient Inception V2 based Deep Convolutional Neural Network for Real-Time Hand Action Recognition. IET Image Processing, 2019, 14(4):688-696.

[13] Carvalho K, VT Basílio, Brando A S. Action recognition for educational proposals applying concepts of Social Assistive Robotics. Cognitive Systems Research, 2022, 71:1-8.

[14] Hu Z, Pan Z, Wang Q, Yu, L, Fei S. Forward-reverse adaptive graph convolutional networks for skeleton-based action recognition. Neurocomputing, 2022, 492(1):624-636.

[15] Gao B K, Dong L, Bi H B, Bi Y Z. Focus on temporal graph convolutional networks with unified attention for skeleton-based action recognition. Applied Intelligence: The International Journal of Artificial Intelligence, Neural Networks, and Complex Problem-Solving Technologies, 2022, 52(5):5608-5616.

[16] Ren Y, Zhang R, Feng Z, Ke C, Yao S, Tang C, Lin L, Ye Y. Macrocephatriolides A and B: Two Guaianolide Trimers fromAinsliaea macrocephalaas PTP1B Inhibitors and Insulin Sensitizers. The Journal of organic chemistry, 2021, 24(21):17782-17789.

[17] J García-Marín, Griera M, R Alajarín, M Rodríguez-Puyol, D Rodríguez-Puyol, Vaquero J. A computer-driven scaffold-hopping approach generating new PTP1B inhibitors from the pyrrolo[1,2-a] quinoxaline core. ChemMedChem, 2021, 18(16):2895-2906.

[18] Chen B, Xingwang J U, Gao Y, Wang J. A Quaternion Two-Stream R-CNN Network for Pixel-Level Color Image Splicing Localization. Chinese Journal of Electronics, 2021, 30(6):1069-1079.

[19] Joardar B K, Doppa J R, Pande P P, Li H, Chakrabarty K. AccuReD: High Accuracy Training of CNNs on ReRAM/GPU Heterogeneous 3D Architecture. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 2020, 5(40):971-984.

[20] Sun K, Zhang J, Liu J, Yu R, Song Z. DRCNN: Dynamic Routing Convolutional Neural Network for Multi-View 3D Object Recognition. IEEE Transactions on Image Processing, 2020, 30(12):868-877.

[21] Yu B, Zhang X, Wu L. A Novel Postprocessing Method for Robust Myoelectric Pattern-Recognition Control Through Movement Pattern Transition Detection.IEEE transactions on human-machine systems, 2020, 50(1):32-41.

[22] Majd M, Safabakhsh R. A motion-aware ConvLSTM network for action recognition. Applied Intelligence, 2019, 49(7):1-7.

[23] Chen H Y, Huang P H, Fu L C. Social crowd navigation of a mobile robot based on human trajectory prediction and hybrid sensing. Autonomous robots, 2023, 47(4):339-351.

# A Lean Service Conceptual Model for Digital Transformation in the Competitive Service Industry

Nur Niswah Hasina Mohammad Amin[1]*, Amelia Natasya Abdul Wahab[2], Nur Fazidah Elias[3], Ruzzakiah Jenal[4],
Muhammad Ihsan Jambak[5]*, Nur Afini Natrah Mohd Ashril[6]

Faculty of Information Science and Technology, Universiti Kebangsaan Malaysia, Selangor, Malaysia[1, 2, 3, 4, 6]
Faculty of Computer Science, Universitas Sriwijaya, Palembang, Indonesia[5]

*Abstract*—In today's competitive service industry, the pressure to boost productivity, cut costs, and improve service quality is immense. By integrating lean principles and digital transformation, organizations can streamline processes and reduce waste. Although various lean models have been developed for different service industry, there is no universal standard. Hence, this study aims to address this gap by proposing a Lean Service Conceptual Model through qualitative research by identifying nine types of waste and seven lean dimensions. Interviews, observations, and audio-visual materials are the data collection methods used in this study. The model aligns seamlessly with modern digital technologies such as big data, the Internet of Things, blockchain, cloud computing, and artificial intelligence, making it adaptable for service organizations to excel in the digital age. The model focuses on enhancing efficiency and effectiveness while primarily reducing waste in service operations. Due to restrictions during the pandemic and the interest expressed by the informants in participating in this study, the focus is thus made on a single case study, which may lead to biased findings. However, future studies will be performed on multiple case studies to enhance the findings. Exploring and reviewing an array of best practices, techniques, and tools available for waste reduction within organizational operations is paramount.

*Keywords*—*Lean principles; digital transformation; model conceptual; service industry; waste; dimension; qualitative research*

## I. INTRODUCTION

In the evolving environment of contemporary companies, two fundamental paradigms have emerged as transformative forces that shape the way organizations operate and deliver value, which is known as lean and digital transformation. The term "lean" refers to the principles and procedures of the Toyota Production System (TPS) [1]. Nowadays, lean is no longer limited to the manufacturing industry; it is all about doing more with less. Over time, the service industry has also effectively adopted the lean idea [2]. Without doubt, lean is successfully implemented in the service industry.

In today's competitive world, service organizations are faced with enormous pressure to raise productivity, cut costs, and enhance the quality of their service. Lean is a practice used in overcoming such issues, and in creating more value for the customers. Studies on the implementation of lean service have indicated interesting trends that began in 2005, with an increase in the number of studies in the service industry such as healthcare, education, public services, hotels, banking, and

information technology [3]. The modern economy is dependent on the service industry, which is very tightly tied to our daily life [4]. Lean implementation in the service industry has difficulty in assessing efficiency because the service industry focuses non-value-added activities on intangible assets [5]. Lean is necessary in improving performance [6]. Thus lean implementation approaches are used in identifying and reducing wastes [7].

In addition, the era of digital transformation has brought about uncommon technological progress and advancement. To keep up with the current opportunities and trends, organizations must therefore constantly innovate through a process known as digital transformation [8]; which is a process of building dynamic capabilities for the continuous strategic regeneration of organizations [9]. Digital transformation is the adoption of digital technology into all aspects of an organization's operations that lead to a significant change in the way the organization operates and gives value for its customers [10]. The approaches towards digital transformation would adopt a different viewpoint that is aimed at accomplishing multiple goals [11]; impacting all aspects of the organizational process, activities, and structures [12].

Despite the growing acknowledgment of the significance of lean principles and digital transformation within the service industry, a discernible gap persists in the existing body of knowledge. While previous studies have delved into either lean practices or digital transformation individually, there remains a shortage of comprehensive research exploring the synergistic effects arising from the integration of lean principles with digital transformation in service organizations [13]. This study endeavours to address this gap by providing a holistic examination of how the convergence of lean principles and digital transformation can contribute to waste reduction in the operational service industry [14].

Our research aims to fill this void by offering valuable insights that not only tackle the current challenges faced by service organizations but also furnish a nuanced understanding of the potential benefits and challenges associated with this integration. Through this endeavour, we seek to provide practical guidance for practitioners, researchers, and policymakers in navigating the intricacies of contemporary service environments, with the goal of facilitating operational efficiency and waste reduction.

The integration of both lean principles and digital transformation has the potential to change the service industry

---

*Corresponding Author.

by enabling organizations to simplify processes, and thus reduce wastes. However, this integration is challenging due to the complexity of the technological requirements, culture, behavioural challenges and also the size of the organization [8]. Over the years, several lean models have been developed for different service industries. However, there is no standardized lean model that can be used for all service industries that have type of waste and dimension all together. Thus, the crucial need for future research to develop a standardized model for lean services [15]. Developing a standardized model is crucial because it ensures uniformity and comparability across studies, allowing for greater validation and reproducibility. Hence, rather than developing a completely new model, this study attempts to examine and employ appropriate models that are currently available for the implementation of lean services.

This paper aims to explore how lean principles and digital transformation may reduce waste in service organizations. We want to propose a conceptual model that enables service organizations to successfully navigate changes in this era of digital transformation, thus maximising their operational efficiency by reducing waste in the operation. Besides that, the scope of this study is to find the answers to the following research question.

RQ: How can the integration of lean principles and digital transformation lead to waste reduction in service organizations?

The organization of this article is as follows. Section II discusses the related works on lean service and data transformation. Section III presents the research method and activities and discusses the approach taken in conducting this study. Section IV presents the results, and Section V presents the discussion. Finally, the conclusion of this study is made in Section VI.

## II. RELATED WORKS

### A. Definition and Principles of Lean

Lean is a philosophy, a collection of lean techniques or tools, and the concept of waste elimination [16]. Lean is the most important word in any organization [17]. The goal of implementing lean in the service industry is the same as the manufacturing industry: to reduce waste and to enhance resource efficiency. Lean defines as the continuous elimination of waste in all areas of operation [18]. To stay ahead of their competitors, the service industry must address the needs of every consumer. Lean service operation must provide what the consumer wants, where he wants it [19].

The manufacturing industry is differentiated from the service industry in terms of waste and dimension. The manufacturing industry is involved in transforming goods or raw materials into new products such as machinery, computers, electronics, furniture, chemicals, food, and plastics [20]. On the other hand, the service industry creates value, particularly intangible values such as management, guidance, information, advice, design, data, and experience.

The two industries have different outputs, demands, customer-specific production, labour requirements, automated processes, and the location of physical production [21]. Table I

illustrates the differences between the manufacturing and the service industry:

TABLE I. THE DIFFERENCE BETWEEN THE MANUFACTURING INDUSTRY AND THE SERVICE INDUSTRY

| Differences | Manufacturing | Service |
|---|---|---|
| Output | Physical products that is observable and touchable by the customer. | Intangible. |
| Demand | Produces product stocks with inventory level that are parallel to the forecast of customer demand. | Does not keep inventories; service is provided as per customer request |
| Customer-specific production | Production can be performed without customer orders or customer demand forecast. | Service is provided only upon customer request. |
| Labour requirements and automated processes | Automating production process to reduce labour needs. | Need to recruit people with specific knowledge and skills. |
| Physical production locations | Must be physically found for production operation and stock keeping. | Does not require physical site for production. |

Source: [21]

The difference in operation makes the service industry unique. It is important in generating economic growth [22]. With a contribution of more than half of a country's gross domestic product, the service industry is indispensable to the global economy [4], [11]. The service industry plays an important role in the global economy [23]. Many services industry try to distinguish themselves from their competitors by making improvements in their operation. Considering the current economic situation, the successful implementation of lean is expected to contribute to cost reduction and improvement of the service operation, [2]; which is achieved through the identification and elimination of waste [24].

Although lean in the service industry began in the 21st century and is continuously expanding [25], it is still a relatively new concept and has not been thoroughly researched [26]. The advantages of lean include reduced inventory, increased process understanding, operational cost reduction, less re-work, reduced lead-time, and less process waste [27]. Although the implementation of lean in the service industry offers numerous advantages, it is nonetheless challenge-free [28]. The main challenge is the lack of awareness regarding its advantages [15].



Fig. 1. Five principles of lean thinking.

Lean Thinking is a principle that could improve the efficiency of the service industry, reduce its operating costs, and increase operations capability [29], [30]]. Over time, Lean Thinking has expanded tremendously in the service industry, providing excellent benefits [31]. Lean becomes a way of

thinking, whereas practices or tools are ways to put these beliefs into reality [15]. Fig. 1 show the five principles of Lean Thinking are; specify what creates value, identify the value stream, flow, pull, and strive for perfection [23].

### B. Models of Lean

This section examines several models of lean based on the review of the literature made between the years 2004 to 2022.

A model for lean production is explained and translated into service industry [32]. The Ahlstrom model is the earliest model designed to assess the level of lean implementation in an organization. The model indicates; (1) lean development, (2) lean procurement, (3) lean manufacturing, (4) lean distribution form, and (5) lean enterprise or organization that is competitive at the global level. However, this model does not indicate the wastes that may exist in an organization.

Vadivel & Sequeira [3] developed a model to investigate lean service activities and their impact on operational performance in the Indian postal service. However, the development of the model is restricted to their consideration for and selection of methods, tools, and techniques that were only put out in the empirical literature review. Other than that, neither the types of waste nor the way that lean should be implemented in an organization's operations are discussed in their model.

In a study by Sreedharan V et al. [33], a focus group and a structured literature review are used in constructing the Green Lean Six Sigma (GLSS) model for the public industry. This model consists of three different stages; procurement, production, and distribution, whereby the flow of activities starting from the procurement to distribution are depicted in the model. Bajjou et al. [34] developed an input-output model for the construction industry. This input-output model consists of three processes; the input, transformation process, and the output, whereby each process has its own principles. However, no mention of the types of waste is made in these models.

Iranmanesh et al. [35] proposed a model to investigate the effects of lean practices involving the aspects of process and equipment, manufacturing planning and control, human resources, product design, supplier relationship, and customer relationship. However, no discussion is made on the effects of waste on the sustainable performance of the manufacturing firms. A product-service system (PSS) leanness assessment model is developed by Elnadi & Shehab [31]. This assessment model consists of three levels; the enablers, criteria, and attributes that are used in proposing an index to assess the leanness of PSS in a United Kingdom manufacturing company. However, this model does not involve the assessment of any type of waste.

In addition to that, in a study conducted by Abdul Wahab et al. [36], a model of lean production dimensions and its relation to waste has been developed. This model serves as a guideline for management team in examining the types and places that waste can occur in the manufacturing industry. The model consists of seven dimensions that are the functional areas in the manufacturing industry. It also states eight types of waste that might exist in each dimension. These seven dimensions are Supplier Relationship, Customer Relationship, Product

Development and Technology, Manufacturing Process and Equipment, Manufacturing Planning and Scheduling, Customer Relationship, and Visual Information System. The eight types of waste are waiting, defect, overproduction, transportation, motion, inventory, extra processing, and underutilized people. However, the dimensions and types of waste presented in this model are specified for the manufacturing industry, hence this model is not applicable for the service industry.

A model is therefore required in providing the right guidance and directions for industries, specifically those in the service industry, to enhance their operations. Although the process of becoming a lean service organization takes time and effort, the development of a model for guiding and tracking the results of such effort is of paramount importance to speed up the process.

### C. Types of Waste

Lean is the outcome of Taichii Ohno's invention of the Toyota Production System that aims to reduce waste. To identify the "*Muda*" or waste in lean service, a review of the literature of previous studies is performed. Waste is defined as any activity that increases the cost, but does not add any value from the customers' perspective [26], [37], [38]. Besides that, lean is about improving quality to eliminate waste [15]. Lean is also an approach of eliminating waste in a process and creates value for the customers [39].

Identifying waste in a service industry can be complex because the operations are intangible [23]. Several types of waste identified by Ohno [40] in the manufacturing industry also exist in the service industry: over-production, inventory, waiting, motion, transportation, defects, and over-processing. Waste in the form of underutilized resources and a manager's resistance to change are also mentioned in the service industry [1], [23].

The discussion on lean service by Mohammad Amin et al. [24] examines nine types of waste; over-production, inventory, waiting, motion, transportation, defect, over-processing, underutilized resources, and manager's resistance to change, which is depicted in Fig. 2.

In this work, the types of waste are identified based on their definition. Table II shown the definition of waste based on the service perspective.



- Over-production
- Inventory
- Waiting
- Motion
- Transportation
- Defects
- Over-processing
- Underutilized resources
- Manager's resistance to change

Fig. 2. Waste in lean service.

TABLE II.    THE DEFINITION OF WASTE FROM THE PERSPECTIVE OF SERVICE INDUSTRY

| Waste(s) | Definition |
|---|---|
| Inventory | Any work in process (Work-in-Progress) that exceeds what needs to be produced for the customer, inventory accumulation that causes the overuse of storage space, and reduced worker productivity due to the surplus of inventory [23], [41]. |
| Transportation | Unnecessary movement of materials, products, information, workers and forklift operators [23], [41]. |
| Waiting | Waiting whereby the employees or customers must wait for information or service delivery. Waiting is also involved when employees are ready to resume work, yet are unable to do so due to product, machine, or system unavailability [1], [23]. |
| Motion | Unnecessary movement of resources or workers that need to bend over to choose items [1], [23]. |
| Over-production | The completion of more work than required or before customer demand, which can lead to overcrowding [1], [23]. |
| Over-processing | Adding unnecessary value to a service or product that is not requested by the customer, or will pay for, including unnecessary inspection and packaging [1], [41]. |
| Defect | Any aspect of the service that does not suit the customer's needs, such as selecting incorrect items or incorrect quantity of an item [1], [23]. |
| Underutilized resources | Waste of resources, especially human potential, not utilizing the talent and potential of employees, underutilizing their skills, creative abilities, and knowledge [1], [23], [42]. |
| Manager's resistance to change | The attitude of "saying no" by the management, does not encourage all employees to be involved in the continuous process of improvement [23]. |

## D. Dimension

The understanding of the dimension of lean implementation is an important aspect in improving operational performance. Lean implementation dimension does not only serve as a strategic guide in identifying and overcoming waste in the service process; it also acts as a guide for the formation of an organizational culture that focuses on efficiency, quality, and customer satisfaction.

Several variations of the functional domains or operational dimensions used in assessing the level of organizational implementation have been identified in the previous studies. However, most studies would focus on the dimensions without discussing the situations where waste could occur. Table III below depicts the definition of the dimensions of lean in the service industry.

## E. Digital Transformation

Digital transformation has become one of the most discussed topics, and many industries have embraced digital transformation to acquire a competitive edge and maintain their sustainability [49]. Service industry operate their businesses, provide customer service and support by changing their way of operating through digital transformation [50].

Amidst the dynamic business landscape changes, service organizations are adopting digital transformation. Digital transformation enables service organizations to manage their operations more efficiently and effectively through the reduction of operational waste. At present, the term digital transformation does not have any recognized definition [51] since the scholarly literature lacks specific definitions [52].

Table IV illustrates some of the definitions of digital transformation [52].

There are new technologies that have become the trend in digital transformation, for example big data [52]–[54], the Internet of Things [9], [52]–[54], blockchain [9], [52], cloud computing [9], [53], [54], and artificial intelligence [52], [54]. These technologies offer new uses based on innovation and focus on the needs of the consumers [52]. Fig. 3 shows the technologies for digital transformation.

TABLE III.    DEFINITION OF LEAN DIMENSION

| Dimension | Definition |
|---|---|
| Lean Supplier | A supplier is a person or company that provides goods or services to another person or entity. The seller is referred to as the supplier. The basic function of supplier management is to control cost, quality, delivery performance, and billing accuracy. Adapted from [43], [44]. |
| Lean Workforce Management | Workforce management is the process of strategically optimizing employee productivity to ensure all resources are in the right place at the right time. Workforce management strategies include scheduling, forecasting, skills management, punctuality and attendance, daily management, and employee empowerment. Adapted from [45]. |
| Lean Operations Development and Technology | Lean Operations Development and Technology refers to the choice of operational structure, materials, and technical solutions in adopting service methods in line with the latest technology or innovative practices to increase operational capability. Technology is used to increase the autonomy of result-oriented groups and the distribution of responsibilities in operations. Adapted from [32], [46]. |
| Lean Service Provision Process | Lean Service Provision Process refers to all activities required in producing services by using the collection of methods and materials or techniques in service operations that emphasize service quality standards, workplace layout, productive use of equipment and maintenance, material handling, safety, hygiene, and ergonomic aspects to reduce service preparation time. Adapted from [46]. |
| Lean Service Planning and Scheduling | Lean Service Planning and Scheduling refers to all activities required to coordinate services and market demand, and thus increase the ability to meet customer orders. This minimizes variation in service operations, which can be achieved by optimizing resource use into a seamless service flow and by maximizing productivity through usage of appropriate service scheduling methods, and tools or techniques. Adapted from [46]. |
| Customer Relationship | Customers are people or organization who receives, use, or purchase products or services, and they can choose different goods and suppliers. Customer relationship refers to the establishment of a relationship with the customers by obtaining information about their needs and wants customers for better understanding of their preferences. This relationship is also important in deciding the value and quality of service from their perspective, and all worthless activities can be targeted for elimination. Adapted from [46]–[48]. |
| Visual Information System | Visual information system refers to an information system that delivers prompt and useful flow of information to relevant decision-makers to obtain quick feedback and corrective actions. This is achieved by using certain visual tools for different purposes in the workplace, such as visual boards, operational status, and performance information that enable the specific personnel to perform tasks appropriately according to company goals. Adapted from [32], [46], [47] |

TABLE IV.    CURRENT DEFINITIONS OF DIGITAL TRANSFORMATION

| Author (s) | Definition |
|---|---|
| Matt et al. (2015) | Digital transformation strategy is a blueprint that supports organization in governing the transformations that arise owing to the integration of digital technologies, as well as in their operations after a transformation. |
| Hess et al. (2016) | Digital transformation is concerned with the changes digital technologies can bring about in a company's business model, which result in changed products or organizational structures or the automation of processes. These changes can be seen in the rising demand for Internet-based media, which has led to changes in the entire business models (for example, in the music industry). |
| Liere-Netheler et al. (2018) | The use of new digital technologies (social media, mobile, analytics, or embedded devices) to enable significant business improvements (such as enhancing customer experience, streamlining operations, or creating new business models). |
| Horlach et al. (2017) | Digital transformation as encompassing the digitization of sales and communication channels, and the digitization of a firm's offerings (products and services), which replaces or augments physical offerings. Furthermore, digital transformation entails tactical and strategic business moves that are triggered by data-driven insights and the launch of digital business models that allow new ways of capturing value. |
| Westerman et al. (2011) Westerman et al. (2014) Karagiannaki et al. (2017) | The use of technology to radically improve performance or reach of enterprises. |

Source: [52]



Fig. 3.    Technology for digital transformation.

*1) Big data.* Big data is a method and technique to retrieve, collect, manage, and analyse large and complex data in which traditional methods of processing data are difficult [52], [55]. The utilization of big data is also on the rise within the waste management and recycling sector [54]. On the other hand, the use of big data requires for careful planning and implementation [56].

*2) Internet of Things (IoT).* This is one of the technologies essential in the evolution of services, and in increasing customer value [57]. IoT involves the connectivity of physical objects to the Internet or other interconnected systems using sensors and actuators [52]. The communication and exchange of data among physical objects can be performed using IoT. The progress in IoT is not restricted solely to Industry 4.0, as it is also concurrent with the evolution of the service transformation [57].

*3) Blockchain.* Currently, several study fields are paying attention to a new technology known as blockchain [58]. which has become a top technology layer for financial applications [59]. It is practical and appropriate for network providers to trade processing and networking resources using a blockchain-based solution [60]. Blockchain has generated interests as an innovative technology that has the potential to provide substantial cost reductions by allowing transactions to be carried out as peer-to-peer operations directly between users [61]. Using blockchain platforms for service institutions is crucial for specific purposes [62].

*4) Cloud computing.* Cloud computing makes it possible for information to be distributed effectively, regardless of the location [53]. This technology plays a significant role in the service industry; customers want to reduce costs, whilst cloud computing service providers provide their customers with services that maximize their earnings [63]. It refers to the delivery of various computing services, and builds on well-established trends for reducing the cost of service delivery [64].

*5) Artificial Intelligence (AI).* Artificial intelligence allows for precise decision-making that offers significant time and cost savings through data collection, forecasting, and trend analysis [65]. For the last two decades, AI has greatly improved the performance of the manufacturing and service industries [66]. AI can also be used for a wide range of tasks, such as identifying data trends to reduce market risks, improving customer service with the help of virtual assistants, and analysing large document repositories spread across numerous servers within an organization to find instances of compliance violations [65].

Digital transformation is reshaping the service industry in profound ways. The integration of technologies such as big data, IoT, blockchain, cloud computing, and AI is driving a fundamental shift in the way service organizations operate and engage with their customers. This transformation enhances operational efficiency, and enables for personalized customer experiences. By using these digital technologies, service providers may go beyond the customers' expectations.

## III.    METHOD

### A. Research Design

During the research design phase, we meticulously reviewed current methodologies employed in studying lean methods within the service industry. Our investigation entailed a comprehensive examination of a variety of qualitative and quantitative approaches. Table V shows the quantitative versus qualitative approaches.

The quantitative approach involves statistical analysis to analyse trends and relationships, comparing results with

predictions and past research, while the qualitative approach focuses on descriptive data analysis, identifying themes through text analysis, and interpreting findings within the study's context [68]. The case study approach encompasses a set of methods that emphasize the choice between a qualitative or quantitative approach [69].

TABLE V. QUANTITATIVE VERSUS QUALITATIVE APPROACHES

| Quantitative approach | Qualitative approach |
|---|---|
| Measure objective facts | Construct social reality, cultural meaning |
| Focus on variable | Focus on interactive processes, events |
| Reliability the key factor | Authenticity the key factor |
| Value free | Values present and explicit |
| Separate theory and data | Theory and data fused |
| Independent of context | Situationally constrained |
| Many cases, subjects | Few cases, subjects |
| Statistical analysis | Thematic analysis |
| Researcher detached | Researcher involved |

Source: [67]

A qualitative method was used in this study to achieve the aim of gaining understanding about lean in the service industry. The decision to adopt a qualitative approach is based on its capability to thoroughly explore the detailed complexities and subjective aspects of implementing lean principles within the service industry.

A qualitative method is an approach that requires the researcher to approach the study subject directly; to observe, listen, ask, and verify [70]. This method explores the informants' perspectives in order to comprehend a group or phenomena [71]. The purpose of qualitative data collection is to determine the types of data that will answer the research questions [68]. As shown in Fig. 4, this study was conducted in three phases, which included (1) Data collection, (2) Analysis, and (3) Result.



Fig. 4. Research design.

### B. Phase 1: Data Collection

Data collection for this study was performed through focus group interview, observation, and audio-visual materials.

*1) Interview protocol:* An interview protocol was developed to validate the conceptual model. Questions on the interview protocol were used to collect data for model validation. The goal of a study is outlined in its interview protocol [72]. The interview protocol used in this study is a semi-structured interview; the interview is performed based on the questions and sequence of questions pre-determined by the

interviewer, and the important content is recorded during the interview session [71].

The interview protocol was constructed by designing questions based on the components in the conceptual model. There are thirty-nine questions that are divided into five parts; Part A: Demographics, Part B: Dimension of Lean in Services, Part C: Relationship among the Dimensions of Lean, Part D: Relationship between the Dimension of Lean and Waste, and Part E: Role of Information Technology in Lean Implementation. Table VI shows the number of questions for each section in the interview protocol for model validation.

The content of the interview protocol has been validated prior to the interview session. The experts in lean commented on every question of the protocol. Four experts took part in validating this protocol; two are academicians actively engaged in lean research, and the other two are from the industry with knowledge of lean. These four experts were approached and invited by e-mail; they were given the protocol interview and one week to complete the evaluation.

A document having the interview questions was given to each informant prior to the interview so that they are familiar with the questions to be posed during the interview session. The interview began with an explanation of the lean conceptual model for the service industry before moving on to the structured questions.

*2) Focus group interview:* The purpose of the focus group interview is to validate the conceptual model. Focus group interview can be used to discuss issues at a more strategic level [73]. Focus group interview allows for multiple informants to be simultaneously interviewed [72]. A purposive sampling procedure was used to choose the sample for this study. However, as purposive sampling is a frequent case study methodology strategy and will yield the most information about the subject under study, snowball sampling was used [74].

Focus group interview usually consists of four to six informants [68]. Thus, four informants were identified from the researcher's initial contact with workers from the case study companies, and were then selected for the interview. A network was later created by asking41 the first group of informants to refer more informants for the focus group interview. Table VII is four informants who expressed their interest in taking part in this study.

TABLE VI. NUMBER OF INTERVIEW PROTOCOL QUESTIONS

| Part | Total |
|---|---|
| Part A: Demographics | 8 |
| Part B: Dimension of Lean in Services | 8 |
| Part C: Relationship Among Dimensions of Lean | 4 |
| Part D: Relationship between Dimension of Lean and Waste | 17 |
| Part E: Role of Information Technology in Lean Implementation | 2 |
| Total question(s) | 39 |

TABLE VII.    LIST OF INFORMANTS FOR INTERVIEW

| Informant number | Role | Years of experience |
|---|---|---|
| IN1 | Warehouse Executive | 16 |
| IN2 | Cargo Operation Executive | 3 |
| IN3 | Cargo Operation Officer | 1 |
| IN4 | Cargo Operation Officer | 1 1/2 |

Focus group interview informants were contacted to obtain their consent to take part in the interview session. A total of four informants were contacted, to whom the interview protocol was sent via email prior to the interview session in preparation for the interview.

The focus group interview session lasted for 73 minutes involving a conversation between the researcher and the informants to obtain relevant data and information. The data and information include; types of dimensions in services, the relationship among the dimensions, the relationship between waste and the dimensions, and the role of information technology in the implementation of lean.

The interview session was physically conducted at the meeting room of the case study company on September 1, 2022. Prior to the session, permission was sought from the informants for the conversation to be recorded [75] using a voice recorder. Although the researcher controlled the discussion by asking questions based on the interview protocol questions, the informants were given the opportunity to speak and share their views freely.

*3) Observation:* The purpose of the observation is to validate the conceptual model. Observation is one of the processes of gathering information openly, directly by observing people and places at the site of the study [68]. The data obtained from the observation in this study are in the form of audio-visual materials and field notes. Since permission to observe the interview was granted for only 60 minutes, it was conducted on September 1, 2022, at the case study company. The researcher was accompanied by three employees from the case study company who understand the operational processes of their company. They consist of a warehouse executive, a cargo operation executive, and a cargo operation officer.

The role of the researcher was only as an observer, not a participant. Non-participant observers are observers who visit the site and record data without being involved in the activities of the participants [68]. Observation of the operational process was performed with the guidance from the three workers who explained the activities that take place in each process.

*4) Audio-visual material:* Data collection was audio-visual materials that made up of photos or sounds of people or locations captured by the researcher or another person to assist the researcher in comprehending the core phenomenon under investigation [68]. During the observation, audio-visual materials such as pictures and videos of the process were taken.

*C. Phase 2: Analysis*

Qualitative data analysis involves the systematic process of identifying meaningful information from the data obtained. In this study, qualitative data was obtained from focus group interview, observation, and audio-visual materials. Fig. 5 is a guideline used in analysing the qualitative data [68]. Prior to the data analysis process, focus group interview data and observations were collected and organized into file folders on the computer. Then the interview data was transcribed by the researcher as data preparation for analysis. The collected data were read repeatedly to comprehend the coding of the data. Encoding the data is conducted for analytical reports.



Fig. 5.   Qualitative data analysis process.

*1) Data transcription:* Transcription is frequently used in qualitative research [76]. Transcribing an interview involves converting audiotape recordings into text data [68]. There are no common guidelines or procedures for transcription [77]. Hence, below are four steps used in transcribing the interview of this study:

- Prepare data during interview by recording the conversation.

- Listen to the recording multiple times when transcribing.

- Identify the informants and label them accordingly while transcribing.

- Use timestamps to show when an informant starts or stops speaking.

*2) Themes and code:* Analysis of the data was conducted using a computer software program for qualitative data, such as interview transcripts and pictures, using Atlas.ti. The data was explored and coded by reading all data collection and then employing the codes. Codes were also collected to create themes that were used as the main findings of the study [68]. Before the data were coded, code themes were decided based on the type of waste and dimension of lean service.

*3) Validation and reliability:* Triangulation is the use of different sources of information to help confirm and improve the clarity or accuracy of research findings [73]. Triangulation is also seen as a qualitative research strategy to test validity through the convergence of information from different sources [78]. The same interview protocol was used for all case study informants during triangulation to increase reliability [79]. This study has chosen a combination of data from focus group interviews and observations to provide triangulation results [80]. Triangulation is used to support the principle in case studies that phenomena are seen and explored from multiple perspectives [81].

*D. Phase 3: Result*

A conceptual model was developed and revised through a case study after the analysis of the data collection was performed. A case study is an intense description and analysis of an experience, social unit, or system related to time or place [82]. A qualitative case study is an ideal method for understanding and interpreting experience. The qualitative case study methodology enables researchers to carry out a thorough investigation of complex phenomena within a particular setting [83].

## IV. RESULTS

To answer the research question "How can the integration of lean principles and digital transformation lead to waste reduction in service organizations?", the qualitative data approach has been carried out. The qualitative method has been successful in gathering feedback from the informants about the dimensions of lean in services, relationship among the dimensions of lean, relationship between the dimension of lean and waste, and the role of information technology in lean implementation. Fig. 6 shows the lean conceptual model that has been developed for the service industry.

In this study, dimensions are defined as functional areas that carry out specific activities and roles in an organization in achieving the organizational goals. There are seven dimensions of lean that exist in the service industry, namely:

- Lean Supplier
- Lean Workforce Management
- Lean Operations and Technology Development
- Lean Service Provision Process
- Lean Service Planning and Scheduling
- Customer Relations
- Visual Information System

Next, in the manufacturing industry, the product production process can be represented as an input-output model, where resources in the form of raw materials will be transformed into finished products due to the output of the system. All informants agreed that all seven dimensions stated can also be represented as input-output processes in the service industry. All the phases in the lean conceptual model will be created based on user or consumer demand.



Fig. 6. Lean service conceptual model (Adapted from [46]).

According to the informants, the lean conceptual model reflects the detailed service operations of the informants' company; business activity must be conducted when there is a demand. All dimensions and their relationships cannot be less than one since the dimensions are interrelated, as specified by the researcher in the initial conceptual model.

Analysis of the data revealed that the informants' organization has nine types of wastes, namely over-production, inventory, waiting, motion, transportation, defect, over-processing, underutilized resources, and manager's resistance to change. According to the informants, the types and examples of wastes are easy to be figured out because they are visible to the naked eyes. Identified waste and its types are particularly important in implement lean service; waste must be identified so that the cause of the problem can be addressed.

TABLE VIII. EXAMPLES OF WASTE IN LEAN DIMENSIONS

| Dimension | Waste(s) |
|---|---|
| Lean Supplier | Defect, Over-production, Waiting, Underutilized resources, Transportation, Inventory, Over-processing. |
| Lean Workforce Management | Defect, Waiting, Underutilized resources, Transportation, Over-processing, Manager's resistance to change. |
| Lean Operations Development and Technology | Defect, Over-production, Waiting, Underutilized resources, Transportation, Inventory, Over-processing. |
| Lean Service Provision Process | Defect, Over-production, Waiting, Underutilized resources, Transportation, Inventory, Over-processing. |
| Lean Service Planning and Scheduling | Defect, Over-production, Waiting, Underutilized resources, Transportation, Inventory, Over-processing. |
| Customer Relationship | Defect, Over-processing. |
| Visual Information System | Defect, Over-production, Waiting, Transportation, Inventory, Over-processing. |

The informants agreed with the proposed dimensions and wastes of the initial conceptual model. However, after analysis, additional waste on several dimensions was discovered as highlighted in the lean conceptual model shown in Fig. 6; over-

production waste is added to the lean supplier dimension, underutilized resources waste is added to the lean operations and technology development dimension, lean service provision process dimension, and lean service planning and scheduling dimension, and defect waste is added in the customer relations dimension. Waste in the seven dimension of the lean conceptual model is illustrated in Table VIII.

According to two informants, IT plays a critical role in assisting with the implementation of lean in the organization. To aid in the application of lean in the service industry, a system must be established. Knowing where waste occurs is a required system feature. However, according one of the informants, they do not require a system to figure out waste because they are more comfortable executing the work manually.

## V. DISCUSSION

This research was carried out to provide a preliminary overview of the validation of a conceptual model for lean service in the service industry. Initially, lean conceptual models are developed based on the literature review and preliminary research where types of waste and dimensions in lean service are identified. However, in this study, we limit the development of our conceptual model by focusing on the types of waste and dimensions in lean service using a qualitative method. The data for this study were gathered via focus group interview, observation, and audio-visual sources [68].

The results of this study revealed that the proposed Lean Conceptual Model for the service industry is applicable. This is because the service operation in the case study company shares the same dimensions of lean, relationships among the dimensions of lean, and relationships between dimensions of lean and waste. Thus, this study has found nine types of waste with seven lean dimensions.

The types of waste identified are over-production, inventory, waiting, motion, transportation, defects, over-processing, underutilized resources, and manager's resistance to change. The seven lean dimensions identified are Lean Supplier, Lean Workforce Management, Lean Operations and Technology Development, Lean Service Provision Process, Lean Service Planning and Scheduling, Customer Relationship, and Visual Information Systems.

The Lean Service Conceptual Model for the service industry can be aligned with the existing digital transformation technologies such as big data, IoT, blockchain, cloud computing, and AI. By harnessing the capabilities of these technologies, organizations will not only embrace lean principles but also propel their service operations into a new era of efficiency and effectiveness [84].

Across all the Lean Service Conceptual Model's dimensions, big data analytics is crucial for reducing waste [85]. Organizations develop the ability to identify and reduce distinct types of waste through the analysis of significant data produced throughout its service operations. As an example, within the Lean Supplier dimension, data analytics can optimize inventory management, leading to a decrease in excess inventory waste. These analytics can help with resource allocation in the context of lean service planning and scheduling [56] and reduce waste [86].

IoT devices emerge as pivotal assets in the alignment of dimensions within the Lean Service Conceptual Model. These devices assume a crucial role in capturing real-time operational data, seamlessly harmonizing with various dimensions of lean service. They play a crucial role in the efficient use of resources, namely taking care of the Lean Workforce Management component. Moreover, IoT devices effectively monitor the intricacies of service provision processes [87], thereby closely aligning with the Lean Service Provision Process dimension. Furthermore, these tools improve customer experience interaction [88] by encouraging mutually beneficial relationship through the Customer Relationship dimension. The result is a decrease in waste brought on by the ability to make informed decisions made possible by these IoT devices.

Blockchain technology serves as a robust pillar in upholding the core principles of lean service, primarily by instilling trust and transparency, with a particular focus on the Lean Supplier dimension. This technology successfully reduces waste in the supply chain by serving as a strong barrier against flaws and dangerous goods [59]. Moreover, blockchain's capabilities extend to the enhancement of transparency in Customer Relationship dimension, where it securely records interactions and transactions. In addition to fostering more trust, this careful documentation also helps to cut down on processing waste. In summary, the integration of blockchain strengthens lean service by promoting waste reduction, transparency, and trust across the service ecosystem.

In the context of lean service, cloud computing appears as a catalyst for facilitating collaboration [89]. Due to its innate abilities, several lean dimensions can be seamlessly coordinated. TCloud computing transforms into an essential channel for the exchange of real-time information by facilitating improved communication and cooperation across multiple functional areas. Through coordinated efforts and real-time information sharing, this collective method enables organizations to jointly detect and manage waste, strengthening the lean service concepts of efficiency and waste reduction.

Within the context of lean service, AI emerges as a powerful force for automation and greater efficiency. AI proves to be a crucial tool for optimizing operations across all dimensions of lean service thanks to its comprehensive range of automation and predictive analytics capabilities. The Lean Operations Development and Technology dimensions are successfully improved because of how well it performs everyday chores. Additionally, AI is crucial to optimizing resource allocation and integrates perfectly with the Lean Service Planning and Scheduling dimension component. Most significantly, AI helps the Customer Relationship Dimension to offer excellent client experiences [90]. In the process, it simultaneously decreases waste by improving overall process effectiveness, reiterating its function as a major enabler of lean service concept.

## VI. CONCLUSION

This research is aimed at providing a conceptual model that enables service organization to successfully navigate changes

in the environment of the digital era while maximising their operational efficiency by reducing waste in their operation. This research contributes to the types of waste in lean service. Nine types of waste have been identified; over-production, inventory, waiting, motion, transportation, defects, over-processing, underutilized resources, and manager's resistance to change. Seven lean dimensions identified are Lean Supplier, Lean Workforce Management, Lean Operations and Technology Development, Lean Service Provision Process, Lean Service Planning and Scheduling, Customer Relationship, and Visual Information Systems.

This study has successfully validated a Lean Service Conceptual Model for the service industry through the qualitative method by identifying nine types of waste and seven lean dimensions. This research is significant because it proves how well this paradigm aligns with modern digital transformation technologies like big data analytics, IoT, blockchain, cloud computing, and AI. These technologies are essential for reducing waste, optimizing resources, encouraging collaboration, and automating all aspects of lean service. This integration highlights the model's adaptability, positioning it as a catalyst for service organizations to thrive in a digitally transformed landscape characterized by enhanced efficiency and effectiveness, with a primary focus on waste reduction within service operations.

Thus, the contribution of this study provides a solid foundation to ensure efficient achievement or performance in the service industry. The constraint of this study is in the limitation of the number of companies for the case study. The pandemic that hit when the study was conducted has caused the ability to interact with various organizations in the service industry to be limited; only one company was ready to take part as a case study company for this study.

In the realm of future research endeavours, it is recommended for the inclusion of case study companies to be expanded, with a deliberate focus on diverse sectors within the service industry. This approach aims to mitigate potential bias in research outcomes and offers a more comprehensive understanding of how the Lean Service Conceptual Model aligns with digital transformation technologies across different service contexts.

Additionally, exploring and reviewing the array of best practices, techniques, and tools available for waste reduction within organizational operations is paramount. Such investigations can unveil effective measures that organizations can readily implement to enhance operational efficiency and minimize waste, contributing to a more sustainable and lean service ecosystem. These research directions hold the potential to further advance our knowledge and practical insights in the pursuit of lean service excellence.

## REFERENCES

[1] M. Escuder, M. Tanco, A. Muñoz-Villamizar, and J. Santos, "Can Lean eliminate waste in urban logistics? A field study," Int. J. Product. Perform. Manag., vol. 71, pp. 558–575, 2020.

[2] W. Chen, "Research and Application of Civil Aviation Ground Service Management based on Lean Management," Atl. Press, vol. 68, pp. 422–427, 2018.

[3] S. M. Vadivel and A. H. Sequeira, "An Operational Performance of Indian Postal Service using Lean Manufacturing Approach – A Conceptual Model," Proc. Int. Conf. Strateg. Volatile Uncertain Environ. Emerg. Mark., no. July, pp. 318–326, 2017.

[4] W. Jiang, P. S. A. Sousa, M. R. A. Moreira, and G. M. Amaro, "Lean direction in literature: a bibliometric approach," Prod. Manuf. Res., vol. 9, no. 1, pp. 241–263, 2021.

[5] E. A. Kotlyarova, K. F. Mekhantseva, L. S. Markin, and M. O. Otrishko, "Application Possibilities and Standardization Features for Lean Methods in Service Industries," IOP Conf. Ser. Earth Environ. Sci., vol. 666, no. 6, 2021.

[6] F. Pakdil, P. Toktaş, K. M. Leonard, and K. M. Leonard, "Validation of qualitative aspects of the Lean Assessment Tool ( LAT )," 2018.

[7] M. Z. Rafique, S. Mumtaz, M. N. A. Rahman, I. A. Mughal, M. A. Khan, and S. M. Haider, "Wastes in lean production systems," Int. J. Innov. Technol. Explor. Eng., vol. 8, no. 8, pp. 1823–1827, 2019.

[8] Z. Van Veldhoven and J. Vanthienen, "Best practices for digital transformation based on a systematic literature review," Digit. Transform. Soc., vol. 2, no. 2, pp. 104–128, 2023.

[9] K. S. R. Warner and M. Wäger, "Building dynamic capabilities for digital transformation: An ongoing process of strategic renewal," Long Range Plann., vol. 52, no. 3, pp. 326–349, 2019.

[10] C. L. Chang, E. Octoyuda, and I. Arisanti, "The Role of Digital Transformation on Strategic Leader: A Systematic Literature Review," ICBIR 2022 - 2022 7th Int. Conf. Bus. Ind. Res. Proc., pp. 289–294, 2022.

[11] C. Matt, T. Hess, and A. Benlian, "Digital Transformation Strategies," Bus. Inf. Syst. Eng., vol. 57, no. 5, pp. 339–343, 2015.

[12] J. Konopik, C. Jahn, T. Schuster, N. Hoßbach, and A. Pflaum, "Mastering the digital transformation through organizational capabilities: A conceptual framework," Digit. Bus., vol. 2, no. 2, 2022.

[13] A. E. Besser Freitag, J. D. C. Santos, and A. D. C. Reis, "Lean Office and digital transformation: a case study in a services company," Brazilian J. Oper. Prod. Manag., vol. 15, no. 4, pp. 588–594, 2018.

[14] K. Ejsmont, B. Gladysz, D. Corti, F. Castaño, W. M. Mohammed, and J. L. Martinez Lastra, "Towards 'Lean Industry 4.0′–Current trends and future perspectives," Cogent Bus. Manag., vol. 7, no. 1, pp. 0–32, 2020.

[15] S. Gupta, M. Sharma, and V. Sunder M, "Lean services: a systematic review," Int. J. Product. Perform. Manag., vol. 65, no. 8, pp. 1025–1056, 2016.

[16] A. N. Abdul Wahab, M. Mukhtar, and R. Sulaiman, "Lean Production System Definition from the Perspective of Malaysian Industry," Asia-Pacific J. Inf. Technol. Multimed., vol. 6, no. 1, pp. 1–11, 2017.

[17] A. Anuar, D. M. Sadek, L. K. Kheng, N. Othman, and N. A. Nordin, "Could A Conceptual Framework of Lean Healthcare, Safety Climate and Operational Performance Achieving Sustainability?," Int. J. Acad. Res. Bus. Soc. Sci., vol. 12, no. 10, 2022.

[18] P. Molina, K. Nuñez, L. Cantú, B. Villarreal, S. Pedro, and G. García, "Routing Lean and Green in UPS," Int. Conf. Ind. Eng. Oper. Manag., no. 2010, pp. 2577–2586, 2014.

[19] H. dos R. Leite and G. E. Vieira, "Lean philosophy and its applications in the service industry: A review of the current knowledge," Production, vol. 25, no. 3, pp. 529–541, 2015.

[20] J. Spacey, "30 Manufacturing term," 2017. [Online]. Available: https://simplicable.com/new/service-industry. [Accessed: 20-Apr-2020].

[21] I. Linton, "Five Differences Between Service and Manufacturing Organizations | Chron.com," 2019. [Online]. Available: https://smallbusiness.chron.com/five-differences-between-service-manufacturing-organizations-19073.html. [Accessed: 20-Apr-2020].

[22] H. S. Abu Hasim, P. B. Tin, and Z. Darawi, "Analisis keperluan tenaga manusia dalam industri Perkhidmatan di Malaysia," in Persidangan Kebangsaan Ekonomi Malaysia ke VII (PERKEM VII), Transformasi Ekonomi Dan Sosial Ke Arah Negara Maju, 2012, vol. 9, no. 1, pp. 993–1000.

[23] E. Andrés-López, I. González-Requena, and A. Sanz-Lobera, "Lean Service: Reassessment of Lean Manufacturing for Service Activities," Procedia Eng., vol. 132, pp. 23–30, 2015.

[24] N. N. H. Mohammad Amin, N. F. Elias, and A. N. Abdul Wahab, "Identifiying Wastes for the Development of Lean Postal Services," in Proceedings of the International Conference on Electrical Engineering and Informatics, 2021.

[25] S. Gupta and M. Sharma, "Empirical analysis of existing lean service frameworks in a developing economy," Int. J. Lean Six Sigma, vol. 9, no. 4, pp. 482–505, 2016.

[26] M. F. Morales-Contreras, M. F. Suárez-Barraza, and M. Leporati, "Identifying Muda in a fast food service process in Spain," Int. J. Qual. Serv. Sci., vol. 12, no. 2, pp. 201–226, 2020.

[27] T. Melton, "The benefits of lean manufacturing: What lean thinking has to offer the process industries," Chem. Eng. Res. Des., vol. 83, no. 6 A, pp. 662–673, 2005.

[28] M. Alsmadi, A. Almani, and R. Jerisat, "A comparative analysis of Lean practices and performance in the UK manufacturing and service sector firms," Total Qual. Manag. Bus. Excell., vol. 23, no. 3–4, pp. 381–396, 2012.

[29] A. Portioli-staudacher and P. Milano, "Lean Implementation in Service Companies," pp. 652–659, 2010.

[30] A. A. A. Mohammad, "Approaching the adoption of lean thinking principles in food operations in hotels in Egypt," Tour. Rev. Int., vol. 21, no. 4, pp. 365–378, 2017.

[31] M. Elnadi and E. Shehab, "Product-service system leanness assessment model: study of a UK manufacturing company," Int. J. Lean Six Sigma, vol. 12, no. 5, pp. 1046–1072, 2021.

[32] P. Ahlstrom, "Lean service operations: Translating lean production principles to service operations," Int. J. Serv. Technol. Manag., vol. 5, no. 5–6, pp. 545–564, 2004.

[33] R. V. Sreedharan, G. Sandhya, and R. Raju, "Development of a Green Lean Six Sigma model for public sectors," Int. J. Lean Six Sigma, vol. 9, no. 2, pp. 238–255, 2018.

[34] M. S. Bajjou, A. Chafi, and A. Ennadi, "Development of a Conceptual Framework of Lean Construction Principles: An Input-Output Model," J. Adv. Manuf. Syst., vol. 18, no. 1, pp. 1–34, 2019.

[35] M. Iranmanesh, S. Zailani, S. S. Hyun, M. H. Ali, and K. Kim, "Impact of lean manufacturing practices on firms' sustainable performance: Lean culture as a moderator," Sustain., vol. 11, no. 4, 2019.

[36] A. N. Abdul Wahab, M. Mukhtar, R. Sulaiman, and K. Shafinah, "Validating the Relationship Between Lean Dimensions and Wastes: A Pilot Study of Malaysian Industries," Int. J. Eng. Sci. Res. Technol., vol. 6, no. 7, pp. 366–375, 2017.

[37] A. Bahaa, Y. Mostafa, and - Mahmoud, "Enhancing Lean Software Development by using Devops Practices," Int. J. Adv. Comput. Sci. Appl., vol. 8, no. 7, pp. 267–277, 2017.

[38] M. K. A. Kiram and M. M. Yusof, "Lean IT transformation plan for information systems development," Int. J. Adv. Comput. Sci. Appl., vol. 11, no. 8, pp. 473–483, 2020.

[39] L. Rexhepi and P. Shrestha, "Lean Service Implementation in Hospital," 2011.

[40] T. Ohno, Toyota Production System: Beyond Large-Scale Production. New York: Productivity Press, 1988.

[41] M. L. George, Lean Six Sigma for Service: How to Use Lean Speed and Six Sigma Quality to Improve Services and Transactions. 2003.

[42] J. A. Douglas, J. Antony, and A. Douglas, "Waste identification and elimination in HEIs: the role of Lean thinking," Int. J. Qual. Reliab. Manag., vol. 32, no. 9, pp. 970–981, 2015.

[43] R. G. Batson, "Supplier Management in Service Industry: What can be Learned from Automotive Manufacturing?," in Intech, vol. 11, no. tourism, 2018, p. 13.

[44] Saloodo, "Who is a Supplier in business? Logistics Terms and Definitions," 2020. [Online]. Available: https://www.saloodo.com/logistics-dictionary/supplier/. [Accessed: 13-May-2022].

[45] Genesys, "What Is Workforce Management?," 2022. [Online]. Available: https://www.genesys.com/definitions/what-is-workforce-management. [Accessed: 13-May-2022].

[46] A. N. Abdul Wahab, "Kerangka Konseptual Aplikasi Audit Kejat Bagi Industri Pembuatan," Universiti Kebangsaan Malaysia, 2017.

[47] A. M. Sánchez and M. P. Pérez, "The use of lean indicators for operations management in services," Int. J. Serv. Technol. Manag., vol. 5, no. 5–6, pp. 465–478, 2004.

[48] M. B. News, "Customer - definition and meaning," 2022. [Online]. Available: https://marketbusinessnews.com/financial-glossary/customer-definition-meaning/. [Accessed: 13-May-2022].

[49] F. E. Ait-Bennacer, A. Aaroud, K. Akodadi, and B. Cherradi, "Adopting a Digital Transformation in Moroccan Research Structure using a Knowledge Management System: Case of a Research Laboratory," Int. J. Adv. Comput. Sci. Appl., vol. 13, no. 9, pp. 375–384, 2022.

[50] M. Jantti and S. Hyvarinen, "Exploring Digital Transformation and Digital Culture in Service Organizations," 2018 15th Int. Conf. Serv. Syst. Serv. Manag. ICSSSM 2018, pp. 1–6, 2018.

[51] D. Schallmo, C. A. Williams, and L. Boardman, "Digital transformation of business models-best practice, enablers, and roadmap," Int. J. Innov. Manag., vol. 21, no. 8, pp. 1–17, 2017.

[52] M.-I. Mahraz, A. Berrado, and L. Benabbou, "A Systematic Literature Review of Digital Platform Business Models," in The International Conference on Industrial Engineering and Operations Management, 2021, vol. 48 LNISO, no. October, pp. 917–931.

[53] T. S. Ilangakoon, S. K. Weerabahu, P. Samaranayake, and R. Wickramarachchi, "Adoption of Industry 4.0 and lean concepts in hospitals for healthcare operational performance improvement," Int. J. Product. Perform. Manag., vol. 71, no. 6, pp. 2188–2213, 2022.

[54] A. K. Feroz, H. Zo, and A. Chiravuri, "Digital transformation and environmental sustainability: A review and research agenda," Sustain., vol. 13, no. 3, pp. 1–20, 2021.

[55] N. Zulkarnain, M. Anshari, and A. Definition, "Big Data : Concept , Applications , & Challenges," no. November, pp. 307–310, 2016.

[56] S. S. Baawi, M. R. Mokhtar, and R. Sulaiman, "Enhancement of text steganography technique using Lempel-Ziv-Welch algorithm and two-letter word technique," Adv. Intell. Syst. Comput., vol. 843, pp. 525–537, 2019.

[57] A. Eigner and C. Stary, "The Role of Internet-of-Things for Service Transformation," SAGE Open, vol. 13, no. 1, pp. 1–21, 2023.

[58] T. Surasak, N. Wattanavichean, C. Preuksakarn, and S. C. H. Huang, "Thai agriculture products traceability system using blockchain and Internet of Things," Int. J. Adv. Comput. Sci. Appl., vol. 10, no. 9, pp. 578–583, 2019.

[59] G. Perboli, S. Musso, and M. Rosano, "Blockchain in Logistics and Supply Chain: A Lean Approach for Designing Real-World Use Cases," IEEE Access, vol. 6, pp. 62018–62028, 2018.

[60] M. Xevgenis, D. G. Kogias, P. Karkazis, H. C. Leligou, and C. Patrikakis, "Application of blockchain technology in dynamic resource management of next generation networks," Inf., vol. 11, no. 12, pp. 1–14, 2020.

[61] Y. Perwej, "Yusuf Perwej. A Pervasive Review of Blockchain Technology and Its Potential Applications," Open Sci. J. Electr. Electron. Eng., vol. 5, no. 4, pp. 30–43, 2018.

[62] N. N. Pokrovskaia, E. A. Rodionova, I. G. Fomina, M. Z. Epshtein, and D. A. Fedorov, "Blockchain and Smart Contracting in the Context of Digital Transformation of Service," Proc. 2022 Conf. Russ. Young Res. Electr. Electron. Eng. ElConRus 2022, pp. 1727–1731, 2022.

[63] E. WEINTRAUB and Y. COHEN, "Cost Optimization of Cloud Computing Services in a Networked Environment," Int. J. Adv. Comput. Sci. Appl., vol. 6, no. 4, pp. 148–157, 2015.

[64] S. Shilpashree, R. R. Patil, and C. Parvathi, "'Cloud computing an overview,'" Int. J. Eng. Technol., vol. 7, no. 4, pp. 2743–2746, 2018.

[65] A. Prasanth, D. J. Vadakkan, P. Surendran, and B. Thomas, "Role of Artificial Intelligence and Business Decision Making," Int. J. Adv. Comput. Sci. Appl., vol. 14, no. 6, pp. 965–969, 2023.

[66] M. Verma, "Artificial intelligence and its scope in different areas with special reference to the field of education," Int. J. Adv. Educ. Res. 5 Int. J. Adv. Educ. Res., vol. 3, pp. 2455–6157, 2018.

[67] W. L. Neuman, Social Research Methods: Qualitative and Quantitative Approaches, 7th editio. Harlow, United Kingdom: Pearson Education Limited, 2013.

[68] J. W. Creswell, Educational Research: Planning, Conducting, and Evaluating Quantitative and Qualitative Research, 4th editio. Boston,MA, United States: Pearson Education (US), 2011.

[69] G. G. Gable, "Integrating Case Study and Survey Research Methods: An Example in Information Systems," Eur. J. Inf. Syst., vol. 3, no. 2, pp. 112–126, 1994.

[70] I. Ahmad and S. Kamarudin, Metodologi Kajian: Pelbagai Gaya Penyelidikan. 2018.

[71] R. Md Ali, F. Mohd Yusof, and F. Shaffie, Pengumpulan Data Kualitatif Dalam Penyelidikan. Kuala Lumpur: Dewan Bahasa dan Pustaka, 2018.

[72] A. Bolderston, "Conducting a research interview," J. Med. Imaging Radiat. Sci., vol. 43, no. 1, pp. 66–76, 2012.

[73] J. Ritchie and J. Lewis, Qualitative Research Practice A Guide for Social Science Students and Researchers, First Edit. SAGE Publications Ltd, 2003.

[74] M. M. Yusof, J. Kuljis, A. Papazafeiropoulou, and L. K. Stergioulas, "An evaluation framework for Health Information Systems: human, organization and technology-fit factors (HOT-fit)," Int. J. Med. Inform., vol. 77, no. 6, pp. 386–398, 2008.

[75] S. R. A. Ibrahim, J. Yahaya, H. Salehudin, and A. Deraman, "The Development of Green Software Process Model A Qualitative Design and Pilot Study," Int. J. Adv. Comput. Sci. Appl., vol. 12, no. 8, pp. 589–598, 2021.

[76] D. G. Oliver, J. M. Serovich, and T. L. Mason, "Constraints and opportunties with interview transcription," Soc. Forces, vol. 84, no. 2, pp. 1273–1289, 2005.

[77] V. Azevedo et al., "Interview transcription: conceptual issues, practical guidelines, and challenges," Rev. Enferm. Ref., vol. 4, no. 14, pp. 159–168, 2017.

[78] N. Carter, D. Bryant-Lukosius, A. Dicenso, J. Blythe, and A. J. Neville, "The use of triangulation in qualitative research," Oncol. Nurs. Forum, vol. 41, no. 5, pp. 545–547, 2014.

[79] N. Nordin and B. M. Deros, "Organisational change framework for lean manufacturing implementation," Int. J. Supply Chain Manag., vol. 6, no. 3, pp. 309–320, 2017.

[80] A. Alkhoraif and P. McLaughlin, "Lean implementation within manufacturing SMEs in Saudi Arabia: Organizational culture aspects," J. King Saud Univ. - Eng. Sci., vol. 30, no. 3, pp. 232–242, 2018.

[81] P. Baxter and S. Jack, "Qualitative Case Study Methodology : Study Design and Implementation for Novice Researchers," vol. 13, no. 4, pp. 544–559, 2008.

[82] L. D. Bloomberg and M. Volpe, Completing Your Qualitative Dissertation: A Road Map from Beginning to End, 4th Editio. Los Angeles, 2008.

[83] Y. Rashid, A. Rashid, M. A. Warraich, S. S. Sabir, and A. Waseem, "Case Study Method: A Step-by-Step Guide for Business Researchers," Int. J. Qual. Methods, vol. 18, pp. 1–13, 2019.

[84] F. D. Cifone, K. Hoberg, M. Holweg, and A. P. Staudacher, "'Lean 4.0': How can digital technologies support lean practices?," Int. J. Prod. Econ., vol. 241, no. 2017, pp. 1–10, 2021.

[85] S. Gupta, S. Modgil, and A. Gunasekaran, "Big data in lean six sigma: a review and further research directions," Int. J. Prod. Res., vol. 58, no. 3, pp. 947–969, 2020.

[86] J. Corbett and C. Chen, "R60. Big Data Efficiency, Information Waste and Lean Big Data Management: Lessons from the Smart Grid Implementation," in CONF-IRM 2015 Proceedings, 2015, vol. 8.

[87] E. Kadiyala, S. Meda, R. Basani, and S. Muthulakshmi, "Global industrial process monitoring through IoT using Raspberry pi," in 2017 International Conference On Nextgen Electronic Technologies: Silicon to Software, ICNETS2 2017, 2017, pp. 260–262.

[88] V. V. Ratna, "Conceptualizing internet of things (IoT) model for improving customer experience in the retail industry," Int. J. Manag., vol. 11, no. 5, pp. 973–981, 2020.

[89] B. J. White, J. A. E. Brown, C. S. Deale, and A. T. Hardin, "Collaboration Using Cloud Computing and Traditional Systems," Issues Inf. Syst., vol. X, no. 2, 2009.

[90] J. Sujata, D. Aniket, and M. Mahasingh, "Artificial intelligence tools for enhancing customer experience," Int. J. Recent Technol. Eng., vol. 8, no. 2 Special Issue 3, pp. 700–706, 2019.

# The Scheme Design of Wearable Sensor for Exercise Habits Based on Random Game

Youqin Huang*, Zhaodi Feng

School of Physical Education, Gannan Normal University, GanZhou, 341000, China

*Abstract*—**The development of random game theory has enabled wearable sensors to obtain actuator evolution in sports exercise, thus the design of user exercise habits during the exercise process has begun to be studied. Conventional devices only focus on automatic adjustment of sports design, with slight shortcomings in personalization. To address this issue, this study added an anchor node localization device to the adaptive search hybrid learning algorithm and analyzed the exercise goals of athletes. At the same time, a semi definite programming method was installed in wearable sensors to achieve the goal of paying attention to the physical condition of athletes. To verify the performance of the fusion device, this study conducted experiments on the Physical dataset and compared it with three models such as Harris Eagle Optimization. The accuracy rates of designing exercise habits schemes for the four devices were 97.4%, 96.5%, 94.7%, and 91.2%, respectively, indicating that the model has the strongest stability. Under the same running time, the energy loss of this model was 0.11kW * h, which performs the best among the four models. When the athletes are different in age, the F1 values of the four devices are 5.9, 4.5, 4.2 and 3.6 respectively. The results indicate that the proposed fusion model has strong robustness and is suitable for designing exercise habit schemes in the evolution of sports exercise actuators.**

*Keywords*—*Random game; adaptive search hybrid learning algorithm; wearable sensors; physical exercise; evolution of actuators; exercise habits; anchor node positioning; semi definite programming method*

## I. INTRODUCTION

The continuous advancement of technology enables wearable sensors to collect the biological information of athletes, providing them with guidance on exercise plans [1-2]. In recent years, sports exercise actuators have gradually integrated electronic technology, providing a more personalized exercise experience. With the increasing emphasis on physical exercise, wearable sensors, as an emerging technology, have been widely used in the field of physical exercise. It can monitor individual exercise data in real-time, such as heart rate and sleep quality, providing users with scientific health guidance and personalized exercise plans. In addition, some sports exercise actuators are also equipped with computer systems, which can help users to better grasp the exercise effect [3]. Conventional sensors are limited to real-time monitoring of body status, and designing the optimal exercise habit plan based on user characteristics still poses challenges [4]. And this method is only suitable for users who exercise in specific areas. For enthusiasts of special terrains, these methods are prone to over generalization during runtime. To expand the search domain of wearable sensors for physical exercise, this study pioneered the construction of a Stochastic

Games (SG) model to simulate the decision-making process of physical exercise personnel in different contexts.

A wearable sensor was designed based on the Adaptive Hybrid Learning of Search (AHLS) algorithm to monitor athletes. To accurately locate it, this study designed Anchor Node Positioning (ANP) technology in wearable sensors and generated a fusion model (AHLSG-PEA). The main content of the study can be divided into six sections. Section I mainly analyze and summarize the application of the current AHLS algorithm. Related works is given in Section II. Section III introduces AHLS and SG into SG and sensors. Section IV conducts simulation experiments on the Physical dataset. Section V delves into results and discussion. And finally, Section VI concludes the paper. The theoretical significance of this study lies in providing a device for designing exercise plans, aimed at helping users achieve a better exercise experience and thus maintain physical health. The theoretical significance of this study lies in providing a device for designing exercise plans, aimed at helping users achieve a better exercise experience and thus maintain physical health.

## II. RELATED WORKS

In the field of scheme design algorithms, research is widely distributed internationally. Nagal et al. designed a system using a hybrid whale grey wolf optimization algorithm. Their algorithm utilized whale parameters to control and balance the randomness of the strategy. For noise in computer signals, they used a controlled search space for filtering. The performance experiment of the algorithm was conducted through signal-to-noise ratio and its average value relationship was evaluated. This algorithm had better strategy planning ability compared to other technologies [5]. Zhou et al. conducted research on the knapsack problem and considered knapsack preferences to extend quadratic multiple knapsacks. They proposed a hybrid evolutionary search algorithm for backpack strategy analysis and generated new offspring solutions based on the crossover operator of the backpack. The experimental analysis of this algorithm utilized adaptive feasible taboo search to improve the offspring solution. This method could accelerate the generation of candidate solutions and streamline their evaluation to propose the best solution [6].

The methods of scheme design are becoming increasingly widespread, and the design of exercise habits is also becoming popular. Meng et al. considered that the sparrow search algorithm is a metaheuristic optimization method, so they applied it to handle multimodal optimization problems. They first introduced chaotic mapping in their research, while also

using adversarial learning methods to increase the diversity of strategies. To verify its effectiveness, they conducted a large number of experiments in the test suite and demonstrated that this method outperforms conventional optimization algorithms in terms of performance [7]. Suresh et al. believed that the recursive whale algorithm has decision-making ability, so they proposed an optimization method for smart grid utilization by reducing production costs through real-time scheduling. The search behavior of this algorithm was experimentally conducted on a work platform through taboo search. The decision capability analysis of this method indicated that the proposed method has less time for scheme design [8]. Yuan led his research group to study the scheme design ability of the population, and adjusted the coordinate system based on the covariance matrix to move the population towards a more favorable direction. To enhance the suggestion ability of covariance, they learned evolutionary algorithms to improve search efficiency. Compared with other algorithms, experiments have shown that this algorithm is a high-quality algorithm [9]. Rajendran et al. found that manual methods have drawbacks in radiation, so they conducted research on computer-aided diagnosis and believed that the most important step is feature selection. Considering that recent algorithms are prone to falling into local optima, they combined grasshoppers with crows to verify that the proposed multi-layer perceptron has strong feature selection ability. This simulation experiment was conducted on MATLAB and compared with many similar algorithms, and its accuracy, sensitivity, and specificity have been proven to be superior to other algorithms [10]. Zhang et al. believed that the evolution of wind speed prediction actuators has a significant impact on decision analysis in the wind power industry, so a hybrid prediction model was developed. In this model, the secondary decomposition technique used wavelet transform. This technology has the ability to adaptively process data, so that the characteristic components of the signal would also be extracted. The experimental verification was conducted using real-life wind turbines and compared with the same prediction algorithm, proving that the model has the smallest statistical error and the strongest adaptability in performance [11].

Numerous experts and scholars have found that research on the application of KM and Long Short Term Memory (LSTM) is very popular, but research on large-scale datasets is still scarce. This study innovatively links the two and holds significant importance in dataset processing.

## III. AN ANP SOLUTION FOR SG IN EXERCISE HABITS

SG involves the uncertainty faced by participants in exercise decision-making, making it difficult to accurately predict determined exercise habits. To assess the evolving exercise habits of physical exercise actuators (PEA), this study combines the AHLS algorithm and uses ANP to select appropriate exercise habits based on the current search state

during user exercise.

### A. Wearable Sensor Combining AHLS Algorithm with SG

The characteristic of SG is that the actions of physical exercise personnel are influenced by random factors, and therefore cannot be determined through simple optimal strategies. So researchers need to weigh different choices based on probability and adopt appropriate strategies to deal with uncertainty. In practical applications, the analysis of SG requires the use of probability theory methods to determine the optimal strategy and final outcome for each athlete, as shown in Formula (1).

$$\begin{cases} F(X) = d\left(X * (\text{Pr}_1 - \text{Pr}_0)\right)/dt \\ \text{Pr}_0 = X^T * X * X_0 * k \end{cases} \quad (1)$$

In Formula (1), $\text{Pr}_1, \text{Pr}_0$ represents the probability of the athlete generating random ideas. $X$ is the final choice of the sports personnel at the time of the event. The time experienced in this process is denoted as $X^T$, and the cooperative effect they have is represented by $X_0$. $k$ is the environmental parameter of the process. In order to help athletes adapt to this trend, this study added the AHLS algorithm to SG, as shown in Formula (2) [12].

$$AH_x = \sum_{i=1}^{I} \sum_{j=1}^{J} A_i^j * H_{ij} \quad (2)$$

In Formula (2) above, the expected result of SG combined with the AHLS algorithm (AHLSG) is represented by $H_{ij}$. $A_i^j$ represents the process network loss during algorithm operation. In the AHLSG algorithm, AHLS uses domain knowledge and empirical rules to guide the search process and quickly find the expected solution of the algorithm [13]. The advantage of the AHLSG algorithm is that it can flexibly select search strategies based on the characteristics of the problem. During the operation of the AHLSG algorithm, the popular domain will also be determined, as shown in Formula (3).

$$d\Delta Ca/di = C_i * \left(\sum_{i=1}^{I} (Ar_i - Ar_0)\right) - \sum_{t=1}^{T} \sum_{z=1}^{Z} a_t a_z \quad (3)$$

In Formula (3), the running parameters of the AHLSG algorithm are represented by $Ar_i, Ar_0$. The running intervals of the algorithm at the current time and the initial time are denoted as $a_t, a_z$. $C_i$ represents the expected value of the AHLSG algorithm for the research input. The operation of this algorithm requires a large amount of motion habit features to establish machine learning models, which will impose a burden on the computational cost of the algorithm [14]. To reduce the running time of the algorithm, a model was established for this study, as shown in Fig. 1.

Fig. 1. Algorithm flow chart combining AHLS and SG.

In Fig. 1, in the motion data processing of the algorithm, the AHLS algorithm represents the motion contour as a level set, which has good robustness. SG smooths the data to remove noise from motion data. The characteristic of SG is high computational efficiency, but its performance is highly dependent on the quality of training data. So this study introduces ANP in SG, combining anchor nodes with the exercise habits of athletes, as shown in Formula (4) below.

$$\begin{cases} An_1 = (An_0 - 1)^2 * (0.5 * \Delta Ca + 0.25) \\ No_2 = (\partial No_1) * (0.1 * N_2 + 0.75 * N_1) \end{cases} \tag{4}$$

In Formula (4) above, the anchor position of the personnel is represented by $An_1$. $An_0$ is the duration of movement at that location. $N_2$ and $N_1$ represent two different methods of motion at different times, and the stability of these two methods is denoted as $No_1$. $No_2$ represents the positioning of the research hypothesis. This study applies ANP to exercise habit judgment and obtains spatial location information of athletes. Based on the results of this data analysis, this study was able to determine the suitable exercise habits as anchor nodes, as shown in Fig. 2.

Fig. 2 is a motion habit determination method based on anchor nodes. When athletes tend to exercise during a specific time period, the machine will choose that time period as the anchor node. Then, based on the selected anchor node, the motion target of the mover on it is set [15].

Based on the exercise habits of the athletes, the machine provides them with personalized exercise advice to help them better achieve movement on anchor nodes. The comparison method between athletes is Formula (5).

$$Co_1 * \int_{-\pi}^{+\pi} \alpha_1 * \beta_1 d\alpha d\beta = Co_2 \int_{-\pi}^{+\pi} \alpha_2 * \beta_2 d\alpha d\beta \tag{5}$$

In equation (5), the exercise method of the athlete before and after the suggestion is denoted as $\alpha_1, \alpha_2$. $\beta_1, \beta_2$ represents the difficulty coefficient when they are implemented. Both are highly sensitive to the quality of the environment they are in, so variables caused by environmental factors are represented by $Co_1, Co_2$. The limitations of this method are the dynamics and persistence of motion on anchor nodes. To address this issue, this study combines ANP and AHLSG algorithms in sensors to design a novel motion habit scheme, as shown in Fig. 3.



Fig. 2. Determine the exercise habit process diagram as the anchor node.



Fig. 3. Exercise habit design scheme combining anchor node positioning and AHLSG algorithm.

In Fig. 3, this method can help researchers judge exercise habits and choose exercise paths based on the behavioral patterns of the athletes. ANP technology provides location information for athletes, and this method can evaluate the quality of movement of athletes. This study combines this method with motion sensor technology, as shown in Formula (6) [16]. This method can not only classify the user's movement behavior, but also analyze the user's posture change information. By evaluating whether the athlete's movements are correct, their exercise methods, such as running and cycling, can be analyzed.

$$\begin{cases} de_1 = de_k / \mu_k \\ \mu_k = \left[ (de - de_1) / (de_k - de_1) \right]^{0.5} \end{cases} \quad (6)$$

In Formula (6), the initial signal of the sensor is denoted as $de_1$, and its total electrical energy intensity during working time is represented by $de$. $de_k, \mu_k$ represents the peak intensity and energy loss of the sensor signal at the current time. This sensor records motion data through devices such as smart wristbands, which can monitor the movement behavior of athletes in real time and provide real-time motion guidance. This study first helps athletes correct bad habits based on their location information, and then makes actual adjustments to the exercise plan based on the exercise time and weather. The degree of adjustment is Formula (7) below.

$$Ha(x) = -de_k * \int_0^1 (\chi_x + \delta_x) d\chi + \mu_k * \int_0^1 (\chi_x + \delta_x) d\delta \quad (7)$$

In Formula (7), $Ha(x)$ is the current habit of the athlete. $\chi_x$ represents the time interval of exercise. The weather conditions during exercise are represented by $\delta_x$. This study combines wearable sensors with the AHLSG algorithm to obtain the movement trajectories of athletes under different exercise habits, thereby conducting in-depth research on individual exercise habits. As the strategies of athletes change, this method can reveal their decision-making process in exercise habits, providing scientific basis for developing personalized exercise habit improvement plans. This method has the ability to improve individual health levels when applied in the field of sports exercise actuator evolution.

### B. Design of AHLSG -based Sensors in the Evolution of Sports Exercise Actuators

PEA evolution refers to the process of improving PEA performance in terms of design and performance. The function of this mechanical device is to assist in basic physical exercise movements [17]. Conventional actuators only have simple adjustment functions and lack personalized adjustments. To improve this point, this study will design sensors based on AHLSG in PEA, combined with the method shown in Formula (8).

$$Act_a = \left\{ Sen_a * \left\| Sen_b - Act_b \right\|_a^b, Sen \in Act \right\} \quad (8)$$

In Formula (8), $Sen_a, Act_a$ represents the initial data of the sensor and actuator. The elements they combine in the two

are represented by $Sen_b, Act_b$. $Sen, Act$ represents the units to which two devices belong. Through this method, sensors and actuators generate a personalized sports exercise program design device (AHLSG-PEA). This study first used sensors to monitor key indicators of exercise personnel, and analyzed the data collected by the sensors, as shown in Fig. 4 [18].



Fig. 4. Workflow of personalized physical exercise program design device AHLSG-PEA.

Fig. 4 shows the process of designing exercise habits for the AHLSG-PEA device. The data of the sensor module is first scheduled using the AHLSG algorithm, and then a motion data model based on motion personnel is established. PEA plays a clustering role in it [19]. The physiological indicators of personnel include blood flow rate and heart rate information. Finally, this study used the proposed algorithm for real-time data analysis, as shown in Formula (9).

$$\begin{cases} \phi_a = \left[ \varepsilon_a * (Sen_a / Act_a) + \gamma_a \right]^{0.5} \\ \varphi_a = (1 - \varepsilon_b) * (\Delta\Omega)^2 \end{cases} \quad (9)$$

In Formula (9), the physiological indicators of the exerciser are denoted as $\varepsilon_a$. $\varepsilon_b$ is the psychological indicator of the same personnel. The coefficient of conversion between the two is represented by $\gamma_a$. $\Delta\Omega$ represents the clustering center of the PEA work process. The analyzed data can be used to construct suitable exercise plans based on the goals of physical exercise. The plan includes action sequences, action specifications, and exercise intensity, and real-time processing of feedback information based on algorithms. Due to the unique characteristics of athletes, it is necessary to automatically adjust the parameters of the exercise model in order to achieve a wide range of exercise effects. The adjustment method follows Formula (10).

$$\eta_1 = arc \min_{\eta_1 \in \eta} \int \int_{\eta_1}^{\eta} \sqrt{\eta_1 - (\eta_0 + \phi_a \forall \varphi_a)^2} \quad (10)$$

In Formula (10), $\eta_1$ represents the current planning method for exercise habits. $\eta_0$ is the optimal design for studying the preset. The AHLSG-PEA device has an evolutionary optimization effect on PEA, while setting a fitness function in the device to continuously optimize the parameters of the actuator to improve its adaptability [20]. The real-time motion data information of athletes can be fed back to the AHLSG-PEA instrument in real time and provide real-time suggestions. This method can help users adjust their strength and rhythm to achieve better exercise results, as shown in Fig. 5.

Fig. 5.   Real-time feedback process diagram of motion data of AHLSG-PEA instrument.

In Fig. 5, sensor technology and real-time feedback mechanism are combined in AHLSG-PEA. This method can track the human body's motion trajectory in real-time based on the motion contour. In sports actuators, this study will use the AHLSG algorithm to provide real-time exercise recommendations. Then, the raw data collected by the sensor is filtered and denoised to improve the stability of the proposed results. In order to update the motion contour model in real-time, this study incorporated the semi definite programming method into the AHLSG-PEA device. The main purpose of this method is to optimize the motion parameters, and the optimization method is Formula (11).

$$Po_i * \left(-1 + \sum_{i=1}^{I} \coprod_{j=1}^{J} Se_i / Se_j \right) = 1 + \sum_{i=1}^{I} \coprod_{j=1}^{J} \left(Se_i + Se_j\right) \quad (11)$$

In Formula (11), the weather conditions during exercise are denoted as $Po_i$. $Se_i, Se_j$ represents the position and mood parameters of athletes in the semi definite programming method. In the optimized AHLSG-PEA device, this study uses rules for information exchange, enabling information to spread across the network. In this device, this propagation method is referred to as rule propagation. The rule propagation model is used to help instruments understand information, where the propagation path is a key factor in the propagation process, represented by Formula (12).

$$Sp(e) = Sp_1 + Sp(e-1)/10 * \left[ \lg\left(e_0 / e\right) \right] \quad (12)$$

In Formula (12), $Sp(e)$ represents the dissemination method of the exercise participants at the current location. The propagation path of the previous athlete is represented by $Sp(e-1)$. $Sp_1$ represents the perfect path preset by the AHLSG-PEA device. The communication barriers between them are denoted as $\lg\left(e_0 / e\right)$. This method can predict the behavioral information of athletes. By simulating their propagation process, this study can evaluate the impact of different strategies on exercise habits and guide the development of exercise methods. In addition, the rule propagation method can also set the optimal motion node, as shown in Fig. 6.

In Fig. 6, when setting the optimal motion node, this study first collects data related to the motion node, including information on population density and distribution of motion facilities. Then, based on the physical exercise goals of the athletes, evaluation indicators are determined to measure the suitability of each potential exercise node. Based on the collected data and determined evaluation indicators, the suitability of each potential motion node is evaluated, as shown in Formula (13). This formula can be used to analyze the evolving motion habits of actuators.

$$Pot_m = \left(Pot_1 / Pot_n + 1\right) \ominus Sp(e) - \left(Pot_0 / Pot_n + 1\right) \quad (13)$$



Fig. 6.   Process diagram of method for setting optimal motion node through rule propagation.

In Formula (13), $Pot_1, Pot_0$ represents the motion nodes in different regions. $Pot_m$ represents the motion node selected as the latent state. The optimal state of this point is represented by $Pot_n$. After obtaining potential motion nodes, suitability scores were assigned to each node. Based on the evaluation score, the best motion node will be selected. Finally, this study suggests that there may be noise issues in real-world nodes, so adjustments are made based on feedback information, as shown in Formula (14).

$$Fee(n) = \sum_{i=1}^{I} \left[ (Fee_i - Fee_1) * (Fee_1 - Fee_m)^2 \right]^{1/2} / \sum_{m=1}^{M} Pot_m \quad (14)$$

In Formula (14), the first feedback information of the real node is denoted as $Fee_1$, and there are a total of $Fee_i$ feedbacks. Their mean is represented by $Fee_m$. $Fee(n)$ represents the output information of the model in the ideal state, in order to achieve the preset constraints of the research. This method can maximize the exercise effect of athletes, and then transform the problem of designing exercise habits into evolutionary optimization of algorithms to improve the adaptability of AHLSG-PEA devices, as shown in Formula (15).

$$\mu_{eq} = \nu(T) - \theta(T) * \vartheta_{eq} * \varpi \lg(\sigma / \sigma_0) \quad (15)$$

In Formula (15), the performance of the AHLSG algorithm caused by time is denoted as $\nu(T)$. The noise changes are represented by $\theta(T)$. $\vartheta_{eq}$ represents the external environmental factor of the AHLSG-PEA device. $\varpi$ is a variable caused by the device's own factors. $\sigma / \sigma_0$ represents the loss coefficient during signal propagation. Through the above scheme design, the AHLSG-PEA device can achieve personalized exercise plans in the evolution of PEA, thereby improving the safety of exercise. Meanwhile, through continuous evolutionary optimization and real-time feedback, the performance of the actuator can be gradually improved, providing a better exercise experience for exercise users. Stochastic game theory is used to analyze the decision-making rules between physical exercisers, and provides the corresponding strategy updating direction for adaptive search hybrid learning algorithm. At the same time, the adaptive search hybrid learning algorithm constantly optimizes the search process and finds the optimal solution of exercise habits in a more efficient way. The combination of these two methods plays an active role in solving the problem of optimal strategy selection in actuator evolution. Anchor node positioning technology can be used to determine the target position. This paper studies the combination of anchor node positioning technology and adaptive search hybrid learning algorithm to make the random game process better adapt to the needs of sports habits in different environments. And by studying the proposed optimization process, the accuracy of positioning is improved.

## IV. EXAMPLE ANALYSIS OF THE FUSION ALGORITHM AHLSG-PEA IN THE DESIGN OF EXERCISE HABIT SCHEMES

This study conducted experiments on the Physical dataset and compared it with three other models to verify the superiority of the AHLSG-PEA device. This dataset contains a total of 857 exercise habits from different regions, including almost all ages of athletes.

### A. Performance Verification of Wearable Sensors for SG

This study was divided into two groups in a 6:12 ratio for the rational utilization of limited data in the Physics dataset, and algorithm learning and experimental verification were conducted on them respectively. Table I shows the equipment screening and parameters used in the experiment.

This study conducted performance validation of the AHLSG-PEA device after setting parameters according to Table I, and compared it with the AHLS algorithm, Harris Hawks Optimization (HHO) algorithm, and Convolutional Neural Network (CNN). The experimental results are shown in Fig. 7.

TABLE I. EQUIPMENT SELECTION AND PARAMETER DETERMINATION IN PERFORMANCE VERIFICATION EXPERIMENT OF DWKM-LSTMIA ALGORITHM

| Equipment selection | Parameter determination |
|---|---|
| Master client | Intel Yeon E8-2079 |
| Language | Easy Chinese |
| Memory of graphics card | 2T*8 |
| Operating system | Windows 7X |
| Age range of athletes | 5-80 |
| Athletes' sports terrain | Plains, mountains, hills and lakes |
| Sensor wearing habit | Shoulder, hand, neck and ankle |
| Time range of movement | Morning, afternoon, evening, early morning |
| Algorithm working time | 15:22:59 |
| Data set | Physic |
| Execution method | Matlab R2147h |



Fig. 7. Image of performance verification experiment results of AHLSG-PEA equipment.

Fig. 7 shows a comparison of the design capabilities of four different systems for exercise habits. As the model runtime increases, the performance ratings of all four algorithms continue to rise, and the AHLSG-PEA device consistently performs the best. At the same time, the debugging levels of AHLSG-PEA, AHLS, HHO, and CNN equipment were 5.8, 4.6, 4.1, and 3.7, respectively. This indicates that the performance of AHLSG-PEA devices is optimal when the operating environment is the same. To verify the robustness of the AHLSG-PEA model, this study conducted experiments on different exercise times and locations, and the experimental results are shown in Fig. 8.

In Fig. 8 (a), when the model is located in the same region, the performance of all four models shows an upward trend as the running time increases. The working characteristic value of AHLSG-PEA equipment is the highest, at 28. The control conditions in Fig. 8 (b) are the same running time but different exercise locations. In Fig. 8 (b), as the exercise location of the participants becomes increasingly difficult, the calculation accuracy of all four models shows a decreasing trend. The highest accuracy values of AHLSG-PEA, AHLS, HHO, and CNN devices are located at 63, with values of 94.5%, 92.7%, 84.1%, and 82.6%, respectively. This indicates that the

AHLSG-PEA device can accurately design the exercise habits of athletes. However, the above experiments can only demonstrate the internal performance of the equipment, and experimental verification is also required for changes in the conditions of the athletes themselves.

### B. Experimental Analysis of AHLSG-PEA Equipment under the Background of PEA Evolution

To conduct experiments on the performance of AHLSG-PEA equipment in the evolution of PEA, this study focused on the different habits of athletes, as shown in Fig. 9.

Fig. 9 shows the AHLSG-PEA performance experiment under different habits of athletes. In Fig. 9 (a), as the number of participants in the exercise gradually increases, the accuracy of scheme design for all four systems shows an upward trend. Among them, the proposed system has the highest calculation accuracy of 99.7%. In Fig. 9 (b), the scheme design accuracy of AHLSG-PEA, AHLS, HHO, and CNN devices are 97.4%, 96.5%, 94.7%, and 91.2%, respectively. This indicates that AHLSG-PEA can adapt to different time habits in exercise program design. The results of the experiment on athletes from different sports locations are shown in Fig. 10.



(a)Relationship between system performance and working hours

(b)Relationship between system performance and workplace

Fig. 8.    Experimental results of robustness of AHLSG-PEA model.



(a)Accuracy during morning exercise

(b)Accuracy during evening exercise

Fig. 9.    Performance experiment of AHLSG-PEA equipment in different habits of athletes.

(a)Performance experiment of athletes
in plain terrain

(b)Performance experiment of athletes
in mountainous terrain

Fig. 10. Scheme design of AHLSG-PEA equipment when athletes move in different places.



(a) Experiment on judging sports plan
of young athletes

(b) Experiment on judging exercise
plan of middle-aged athletes

Fig. 11. Experimental results of extensive verification of AHLSG-PEA model.

In Fig. 10 (a), the energy losses of AHLSG-PEA, AHLS, HHO, and CNN equipment in plain terrain are 0.11, 0.19, 0.24, and 0.28 kW * h, respectively. This indicates that the proposed model has the highest economic benefits. In Fig. 10 (b), the accuracy of determining the motion habits of the four models in hilly terrain is directly proportional to the running time of the models, and the research model always has the highest accuracy, at 95.8%. This indicates that the AHLSG-PEA model has the best judgment performance. To verify the universality of the model, the results of experiments conducted on users of different ages are shown in Fig. 11.

Fig. 11 shows the extensive validation of the AHLSG-PEA device for young and middle-aged athletes participating in sports activities. In Fig. 11 (a), the F1 value of the

AHLSG-PEA device is the highest, at 5.9, indicating that the model has the strongest stability in experiments with young people. In Fig. 11 (b), the accuracy of AHLSG-PEA, AHLS, HHO, and CNN instruments are 99.8%, 97.2%, 94.1%, and 92.5%, respectively. Therefore, the AHLSG-PEA device has the best performance in designing exercise habits in the evolution of PEA, which is suitable for optimizing the user's exercise experience.

## V. RESULTS AND DISCUSSION

Random game theory is a mathematical model used to analyze the game process of sports. Adaptive search hybrid learning algorithm combines adaptive search and hybrid learning process, while anchor node positioning is used to locate the target position. In this study, these three methods

are integrated to construct a wearable sensor under the background of the evolution of physical exercise actuators, and its application ability in the design of exercise habits is discussed.

The research experiment is divided into two parts, including the test of internal performance and the proof of robustness, effectiveness and universality. The experiment of this model is carried out in the Physic data set and compared with the other four methods. When the running time of the model increases, the performance scores of the four algorithms are rising, and the AHLSG-PEA device always performs best, which shows that the performance of the device is the best. The highest accuracy of AHLSG-PEA device is 94.5%, and its performance is the best among the four devices, which shows that it can accurately predict the exercise habits of athletes.

When the four kinds of equipment are faced with the choice of exercise habits at different times, the accuracy of their scheme design is 97.4%, 96.5%, 94.7% and 91.2% respectively, which shows that the AHLSG-PEA model proposed in this study has good adaptability to different times. The energy loss of AHLSG-PEA instrument is 0.11 kW*h for detecting the exercise habits of athletes in different environments. Compared with AHLS, HHO and CNN, this figure can show the highest economic benefit and the best performance in energy loss. When faced with the detection of exercise habits of people of different ages, the proposed device has the highest F1 value, which shows its strong robustness.

To sum up, the integration technology of random game theory, adaptive search hybrid learning algorithm and anchor node positioning (AHLSG-PEA) is robust, effective and extensive in the design of exercise habit scheme. The experimental results show that the AHLSG-PEA instrument proposed in this study can effectively design the exercise habits of physical exercisers on wearable sensors, and is suitable for being widely used in sports navigation systems. When athletes use non-contact sensors, the research method has certain limitations. With the continuous development of artificial intelligence technology, this integration method will also be more widely studied.

## VI. CONCLUSION

With the evolution of PEA, the recommendation of exercise habits is gradually becoming more personalized. This study added SG technology to the AHLS algorithm and designed it simultaneously with ANP in the sensor, generating the AHLSG-PEA model. To verify its practicality and universality, this study conducted experiments on the Physical dataset and compared the results with AHLS, HHO, and CNN devices. In the internal performance test, the debugging levels of the four instruments were 5.8, 4.6, 4.1, and 3.7, respectively, indicating that the AHLSG-PEA equipment has the best performance. The highest accuracy values of AHLSG-PEA, AHLS, HHO, and CNN devices were 94.5%, 92.7%, 84.1%, and 82.6%, respectively, indicating that AHLSG-PEA devices can accurately predict the exercise habits of athletes. The accuracy rates of scheme design for AHLSG-PEA, AHLS, HHO, and CNN devices were 97.4%, 96.5%, 94.7%, and

91.2%, respectively, for model performance at different times, indicating that the model has good adaptability to different times. For different terrains of athletes, the energy consumption of the four instruments was 0.11, 0.19, 0.24, and 0.28 kW*h, respectively, indicating that the proposed model has the highest economic benefits. When the age of athletes was different, the F1 values of AHLSG-PEA, AHLS, HHO, and CNN devices were 5.9, 4.5, 4.2, and 3.6, respectively, indicating that the AHLSG-PEA device has robustness in the age experiment of athletes. The experimental data showed that the device proposed in this study can effectively cope with internal parameters and external environment, thereby providing users with effective motion plans. However, this study only focuses on wearable sensors, and this method has certain limitations when non-contact sensors are used by athletes. This is because the total amount of data in the dataset is limited. As the number of volunteers increases, this limitation will gradually improve in future research.

## REFERENCES

[1] Wei Q, Chen X. Average stochastic games for continuous-time jump processes. Operations Research Letters, 2021, 49(1): 84-90.

[2] Nsugbe E. Toward a Self-Supervised Architecture for Semen Quality Prediction Using Environmental and Lifestyle Factors. Artificial Intelligence and Applications, 2023, 1(1): 35-42.

[3] O'Brien M K, Botonis O K, Larkin E, Carpenter J, Jayaraman A. Advanced machine learning tools to monitor biomarkers of dysphagia: a wearable sensor proof-of-concept study. Digital Biomarkers, 2021, 5(2): 167-175.

[4] Reynolds M. App Keeps Athletes Hydrated via Wearable Sensor, Smart Bottle, & Refill Pod System.Packaging world, 2022, 29(11):48-50.

[5] Nagal R, Kumar P, Bansal P. A Hybrid Optimization Method OWGWA for EEG/ERP Adaptive Noise Canceller With Controlled Search Space. International Journal of Swarm Intelligence Research (IJSIR), 2020, 11(3): 30-48.

[6] Zhou Q, Hao J K, Wu Q. A hybrid evolutionary search for the generalized quadratic multiple knapsack problem. European Journal of Operational Research, 2022, 296(3): 788-803.

[7] Meng K, Chen C, Xin B. MSSSA: A multi-strategy enhanced sparrow search algorithm for global optimization. Frontiers of Information Technology & Electronic Engineering, 2022, 23(12): 1828-1847.

[8] Suresh M, Meenakumari R. Optimum utilization of grid connected hybrid renewable energy sources using hybrid algorithm. Transactions of the Institute of Measurement and Control, 2021, 43(1): 21-33.

[9] Yuan S, Feng Q. Covariance Matrix Learning Differential Evolution Algorithm Based on Correlation. International Journal of Intelligence Science, 2020, 11(1): 17-30.

[10] Rajendran R, Balasubramaniam S, Ravi V. Hybrid optimization algorithm based feature selection for mammogram images and detecting the breast mass using multilayer perceptron classifier. Computational Intelligence, 2022, 38(4): 1559-1593.

[11] Zhang Y, Zhang W, Guo Z, Zhang J. An effective wind speed prediction model combining secondary decomposition and regularised extreme learning machine optimised by cuckoo search algorithm. Wind Energy, 2022, 25(8): 1406-1433.

[12] Ge J, Liu X, Liang G. Research on Vehicle Routing Problem with Soft Time Windows Based on Hybrid Tabu Search and Scatter Search Algorithm. Computers, Materials & Continua, 2020, 64(3):1945-1958.

[13] Nemmich M A, Debbat F, Slimane M. Hybrid bees approach based on improved search sites selection by firefly algorithm for solving complex continuous functions. Informatica, 2020, 44(2):183-198.

[14] Han P, Guo Y, Li C, H Zhi, Y Lv. Multiple GEO satellites on-orbit repairing mission planning using large neighborhood search-adaptive genetic algorithm. Advances in Space Research, 2022, 70(2): 286-302.

[15] Amin S N, Shivakumara P, Jun T X, Chong L, Zan D L L, Rahavendra R.

An Augmented Reality-Based Approach for Designing Interactive Food Menu of Restaurant Using Android. Artificial Intelligence and Applications, 2023, 1(1): 26-34.

[16] Jia S, Yang C, Chen X. Intelligent Three-dimensional Node Localization Algorithm Using Dynamic Path Planning. Recent Advances in Electrical & Electronic Engineering (Formerly Recent Patents on Electrical & Electronic Engineering), 2021, 14(5): 586-596.

[17] Tian C, Hu X. Mathematical modeling of security impact analysis of communication network based on Monte Carlo algorithm. Computer Communications, 2020, 157(May): 20-27.

[18] Ho Y H, Chan H C B. Decentralized adaptive indoor positioning protocol using Bluetooth Low Energy. Computer Communications, 2020, (Jun.)159: 231-244.

[19] Charles D. The Lead-Lag Relationship Between International Food Prices, Freight Rates, and Trinidad and Tobago's Food Inflation: A Support Vector Regression Analysis. Green and Low-Carbon Economy, 2023, 1(2):94-103.

[20] Zhang L. A BSDE approach to stochastic differential games involving impulse controls and HJBI equation. Journal of Systems Science and Complexity, 2022, 35(3): 766-801.

# The Construction and Application of Library Intelligent Acquisition Decision Model Based on Decision Tree Algorithm

Hong Pan

Library, Yanbian University, Yanji, 133000, China

*Abstract*—In today's digital age, libraries, as the core institutions of knowledge management and information services, are facing an increasing demand from readers. In order to provide more efficient, accurate, and personalized interview services, intelligent interview decision-making in libraries has become an important research field. Traditional manual interview services face challenges such as personnel training and knowledge updates, making it difficult to quickly adapt to new needs and changes. To address these issues, research is being conducted on using machine learning technology to perform post pruning on the basis of standard decision trees and combining it with fuzzy logic to design a fuzzy decision tree. The experimental results show that the F rejection rate (FN) of the model rapidly decreases to about 0.1 as the number of training iterations gradually increases, and stabilizes at around 0.05 after 210 rounds of training, which is 0.10 lower than the rule-based decision model FN. The intelligent acquisition decision-making model designed in this study has higher accuracy and stability, and has certain application potential in the field of intelligent acquisition decision-making in libraries.

*Keywords—Decision tree; machine learning; fuzzy logic; intelligent interview model; post-pruning*

## I. INTRODUCTION

With the continuous development of information technology and the rise of intelligent applications, library management is facing the demand for greater efficiency, convenience, and intelligence. Among them, intelligent acquisition decision-making in libraries is an important link, which involves analyzing, answering, and recommending services to readers. In order to improve the accuracy and efficiency of intelligent acquisition decision-making in libraries, many researchers use collaborative filtering algorithms, which predict books that users may be interested in by analyzing their historical behavioral data. The disadvantage is that it requires a large amount of user behavior data, and the algorithm's recommendation efficiency is not high when new books or new users join (cold start problem) . Some scholars also use hybrid recommendation systems that combine collaborative filtering, content recommendation, and other recommendation methods to compensate for the shortcomings of a single algorithm. However, hybrid recommendation systems may become complex and difficult to manage, and parameter tuning is a challenge [1]. The decision tree model has higher interpretability compared to other machine learning models, such as neural networks or support vector machines [2]. They make decisions through a series of easily understandable rules, enabling library staff to understand the recommendation logic of the model. It achieves data classification and prediction by dividing the dataset into different subsets and continuing to recursively construct decision rules on each subset [3]. In intelligent interview decision-making in libraries, decision tree algorithms can construct corresponding decision rules based on the characteristics of the questions raised by readers, helping the library system respond quickly and accurately to reader needs. However, standard decision trees have defects such as poor processing of continuous data, sensitivity to changes in input data, and susceptibility to overfitting. To address these issues, a fuzzy decision tree (FDT) is designed based on standard decision trees and combined with fuzzy logic, and applied to intelligent acquisition decision-making in libraries. It is expected that optimizing the existing decision tree model will help improve the quality and accuracy of library system interview decision-making. The article mainly consists of four parts. The second part is a review of the current research status on decision trees and fuzzy logic both domestically and internationally. The third part establishes an intelligent acquisition decision model for libraries based on improved decision trees. The fourth part conducts comparative experiments and applicability experiments on the optimization effect of the model.

The novelty of the article lies in the following points. First, the model considers many factors, such as book value, purchase funds, readers' needs and collection structure, and makes acquisition decisions more consistent with the actual needs of library services by building a more comprehensive decision-making framework. Secondly, using information gain as the criterion for feature selection ensures that the model focuses on variables that have a significant impact on classification, such as author, category, price, publisher, and publication time, thereby improving the prediction accuracy of the model. Thirdly, the C4.5 method in the decision tree algorithm was used and pruned to avoid overfitting, thereby optimizing the model's generalization ability. Fourthly, by using fuzzy theory to deal with uncertainty and ambiguity problems, the traditional decision tree algorithm has been enhanced with the ability to handle fuzzy classification, enabling the model to more finely depict real-world situations, especially in situations where user needs are fuzzy or library resource descriptions are unclear. Fifthly, in the process of constructing a decision tree model, by designing fuzzy sets and membership functions, combined with fuzzy logic, the

model can better handle the uncertainty in classification.

This article demonstrates the broad application prospects of artificial intelligence and machine learning technology in library work through the construction and application research of a library intelligent acquisition decision model based on decision tree algorithm. In the future, with the continuous progress of technology and the accumulation of data, intelligent interview decision-making models will play an increasingly important role in the actual work of libraries. Meanwhile, the research methods and achievements of this article can also provide reference and inspiration for intelligent decision-making in other fields.

## II. RELATED WORKS

Recently, with the widespread application of machine learning in various industries, more and more people have begun to use machine learning to solve various difficulties in social learning. Charbuty et al. designed a decision tree based author information classifier to address the issue of low accuracy in traditional author topic classifiers. The experiment showcased that the classification accuracy of the model reached 98.65% [4]. Aldino and other researchers designed a classifier using the decision tree C4.5 algorithm to make it easier for management to determine who is the appropriate student to receive financial aid, and conducted tenfold cross validation on the classification results; The experiment showcases that the accuracy, precision, and recall of the model are all 87%, which means that the model can be well implemented in the system [5]. Tangirala et al. identified a mixture of attributes and minimum class labels for splitting conditions on each non leaf node of a decision tree to solve the problem of homogeneous fruit subsets, and proposed several splitting indices to evaluate splitting; The experiment demonstrates that applying two different segmentation indices, GINI index and information gain, gives the same accuracy [6]. Li and other researchers studied a practical federated environment with relaxed privacy constraints to address the issue of insufficient efficiency or effectiveness of gradient based decision trees for practical applications; the experiment showcases that compared to normal training using local data from each party, this method can significantly improve prediction accuracy [7]. Nancy and others proposed a new intrusion detection system to address the issue of intrusion detection systems neglecting the identification of new types of attacks; the system uses intelligent decision tree classification algorithms to detect known and unknown types of attacks; the experiment demonstrates that this method reduces false positives, energy consumption, and latency [8]. Ramya and other researchers have designed a detection and classification system that combines decision trees with intrusion detection systems (IDS) to enhance the energy security performance of the power grid; The experiment illustrates that the model can accurately classify and predict most attacks [9].

Elhazmi et al. designed a prediction model for mortality of severe adult COVID-19 patients admitted to ICU by combining decision tree and conventional model logic regression to predict the mortality of COVID-19 patients; the experiment showcases that the accuracy of the prediction model reaches 96.65% [10]. Mariniello and other researchers

designed a method category on the ground of decision tree ensemble vibration to address the challenges of structural damage detection and localization in vibration data, and analyzed the dynamic characteristics of the structural system (i.e. pattern shape and natural frequency) to obtain a structural health assessment model; The experiment indicates that the accuracy, reliability of probability prediction, and positioning error of the model perform best when compared with multiple algorithms [11]. Li et al. proposed adaptive control of the gradient of training data for each iteration and leaf node pruning to improve the accuracy of the GBDT model while preserving strong guarantees of differential privacy, to tighten the sensitivity limit; the experiment demonstrates that this method can achieve better model accuracy than other baselines [12]. Sharma and other researchers used fuzzy logic methods to control the operation of ventilation systems that provide fresh air to the environment to solve the problem of insufficient ventilation in indoor environments that damages human health; The experiment showcases that the indoor ventilation system controlled by this model increases the ventilation rate by 15.6% [13]. Arji and others have designed a rule-based fuzzy logic, adaptive neuro fuzzy inference system (ANFIS) to address the significant impact of the spread of infectious diseases on global health and economy; the experiment illustrates that this technology has improved the accuracy of infectious disease identification by 31.34% [14].

In summary, decision trees and fuzzy logic play an increasingly important role in the development of machine learning, but there are few related studies that combine the two to assist in intelligent library acquisition decision-making. Therefore, this study combines fuzzy logic to design a library intelligent acquisition decision model on the ground of decision trees, to further improve the efficiency of library management.

## III. IMPROVEMENT OF DECISION TREE AND DESIGN OF LIBRARY INTELLIGENT ACQUISITION DECISION MODEL SCHEME

This chapter is separated into two sections. The first section first designs a library intelligent acquisition decision model scheme on the ground of the standard decision tree. The second section mainly addresses the shortcomings in the standard decision tree and combines it with fuzzy logic to make some improvements, designing a fuzzy decision tree.

### A. Optimization Strategy for Book Interview Based on Decision Tree

To reasonably select interviewees and provide books and materials that meets the needs of readers in the limited resources of the library. This requires effective interview decision-making methods to improve service quality and promote the dissemination of knowledge and the development of academic research. At present, there are four main influencing factors involved in the decision-making of book acquisition in universities, including book value, procurement funds, reader needs, and collection structure. The collection structure of a library mainly affects the service level and overall level of literature resources. In different application fields, the acquisition concept of libraries also varies, resulting in similar collection structures. The structural framework of

influencing factors for book acquisition decision-making is shown in Fig. 1.

To ensure the effectiveness of the decision model, it is also necessary to select more suitable feature variables, and the selection criteria mainly depend on the classification effect of the samples. The methods for selecting feature variables include information gain or information gain ratio, and this study mainly uses information gain to select feature variables. When learning decision trees, the five features with greater information gain selected include author, category, price, publisher, and publication time. Among them, the author's writing level determines the overall quality of the book content. The author of a book plays an important role in the

procurement of books. Decision tree is a commonly used machine learning algorithm used to solve classification and regression problems. It is a tree based model that performs prediction by segmenting and judging the dataset [15-16]. The decision tree is composed of nodes and edges, with each node representing a feature or attribute, and edges representing the relationship between feature or attribute values [17]. The root node represents the most important feature, the internal node represents the intermediate feature, and the leaf node represents the final output or decision result. The advantages of decision trees include ease of understanding and interpretation, ability to handle discrete and continuous features, and robustness to outliers and missing data. The structural diagram of the decision tree is shown in Fig. 2.



Fig. 1.   The structural framework of factors influencing book acquisition decisions.



Fig. 2.   Structure diagram of decision tree.

Fig. 2 shows a schematic diagram of the decision tree structure, which includes decision binding points, solution branches, probability branches, probability bifurcation points, profit and loss values, etc. Among them, decision junction points, also known as internal nodes or split nodes, represent the decision-making basis for dividing the current data. The decision binding point uses a certain feature or attribute and corresponding threshold to divide the dataset into two or more subsets. The solution branch refers to the edge of the decision junction point, which represents different decision paths or options. Each scheme branch corresponds to a specific feature value or attributes value, indicating that the dataset moves towards different sub nodes on the ground of the value of that feature value. The probability branch indicates the probability transition from one node to another, corresponding to the conditional probabilities of different categories. Profit and loss value is an indicator that measures the effectiveness of decision tree splitting. When selecting a certain feature as the decision binding point, the quality of partitioning is evaluated by calculating the profit and loss value after partitioning on that feature. The profit and loss value can be on the ground of different criteria, such as the Gini index or information gain. The Gini index in the decision tree is shown in Eq. (1).

$$Gini(p) = 1 - \sum (pi)^2 \qquad (1)$$

In Eq. (1), $pi$ represents the probability that the sample belongs to the $i$-th category. The formula for information gain in the decision tree is shown in Eq. (2).

$$IG(D, A) = H(D) - \sum (|\frac{Dv}{D}|) * H(Dv) \qquad (2)$$

In Eq. (2), $D$ represents the dataset; $A$ represents a feature; $H(D)$ represents the entropy of the dataset; $H(Dv)$ represents the entropy of a subset. The information gain in C4.5 algorithm is shown in Eq. (3).

$$GainRatio(D, A) = IG(D, A) / SplitInfo(D, A) \qquad (3)$$

In Eq. (3), $SplitInfo(D, A) = -\sum (\frac{|Dv|}{|D|}) * \log^2 (\frac{|Dv|}{|D|})$ .

Gini index and information gain are commonly used indicators to select the best features for node partitioning. The information gain ratio is an improved indicator introduced in the C4.5 algorithm, which avoids excessive preference for features with more values. Then, it is necessary to prune the decision tree, which is a technique used to reduce the complexity of the decision tree model. When constructing a decision tree, overfitting often occurs, where the model is too complex to generalize well to new data samples. Pruning is the process of reducing the risk of overfitting by pruning some branches or leaf nodes of a decision tree, thereby improving the generalization ability of the model. Decision tree pruning can be divided into two methods: pre-pruning and post-pruning. This study used post pruning, as shown in Fig. 3.



Fig. 3. Decision tree pruning process.

Fig. 3 shows the pruning process of a decision tree, which involves pruning an existing decision tree from the bottom up after it, is constructed. Firstly, this study evaluates each leaf node and calculates its performance indicators (such as accuracy, error rate, etc.) on the validation set. Then, it gradually tries pruning, replacing the leaf node with its parent node, and observes whether the model performance has been improved. If the performance of the model improves after pruning, perform pruning operations; otherwise, keep it as is.

*B. An Intelligent Acquisition Optimization Model for Library Based on Fuzzy Decision Tree*

On the ground of the basic principle of decision tree, establish a library intelligent acquisition decision model on the ground of decision tree. Firstly, it collects data related to library interviews, including user information, library resources, borrowing history, etc. Then there is data preprocessing, which involves cleaning and preprocessing the data, including handling missing values, outliers, and

duplicate values. The third step is feature selection, which evaluates the importance of features through feature analysis and correlation detection. The fourth step is data partitioning, which divides the dataset into training and testing sets for model training and evaluation. Then it constructs a decision tree and uses the decision tree algorithm C4.5 to construct a decision tree model. It then prunes the decision tree and performs pruning operations on the constructed decision tree to avoid overfitting. The seventh step is model training and evaluation, evaluating the performance of the model on the ground of the performance of the test set. Then it uses the constructed decision tree model to predict and make decisions on new interview situations or user needs. Finally, on the ground of feedback and results from practical applications, the model is optimized and adjusted. The basic algorithm flow of the library intelligent acquisition decision model on the

ground of decision tree algorithm is shown in Fig. 4.

Decision tree algorithms are usually able to generate models with good interpretability and comprehensibility, but when the problem has ambiguity, traditional decision trees may become complex and difficult to understand. Library interviews involve many ambiguous situations, such as the ambiguity of user needs and the fuzzy description of library resources [18-19]. Traditional decision trees can only handle discrete classification and attribute values, and cannot effectively handle ambiguity. To address this issue, this study introduces fuzzy theory to handle the uncertainty of data. It allows node partitioning to have fuzzy membership, rather than strict binary partitioning [20]. Such fuzzy partitioning can better handle data with uncertainty or fuzziness. Fuzzy theory can be roughly divided into the following types, as shown in Fig. 5.



Fig. 4.    The algorithm process of library intelligent acquisition decision model.



Fig. 5.    Fuzzy theory classification.

As shown in Fig. 5, fuzzy theory is a very large concept with extensive applications in modern society. It mainly covers aspects such as fuzzy mathematics, fuzzy decision-making, fuzzy systems, uncertainty and information, as well as fuzzy logic and artificial intelligence. Fuzzy mathematics is the foundation of fuzzy theory, which is used to deal with problems that cannot be clearly classified as true or false. Fuzzy decision-making is dedicated to dealing with decision-making problems containing fuzzy factors. Fuzzy system is a method of applying fuzzy theory, which simulates human thinking patterns and generates fuzzy outputs by applying fuzzy rules to input data. Fuzzy logic and artificial intelligence are important components of fuzzy theory, which includes approximate reasoning and fuzzy expert systems. The application of fuzzy theory can better handle uncertainty and fuzziness, thus achieving more accurate and reliable results. Assuming X is a given finite set, the fuzzy subset is shown in Eq. (4).

$$A = \frac{A(e_1)}{e_1} + \frac{A(e_2)}{e_2} + ... + \frac{A(e_N)}{e_N} \tag{4}$$

In Eq. (4), the potential in the set is used to measure the size of the set, as defined in Eq. (5).

$$M(A) = \sum_{i=1}^{N} A(e_i) \tag{5}$$

In Eq. (5), $A$ represents a fuzzy subset. The probability formula for samples belonging to a certain class in nodes in a fuzzy decision tree is shown in Eq. (6).

$$\beta_A^C = SIM(A,C) = \frac{M(A \cap C)}{M(A)} = \frac{\sum_{x \in X} \min(\mu_A(x), \mu_C(x))}{\sum_{x \in X} \mu_A(x)} \tag{6}$$

In Eq. (6), $M(A)$ represents the sum of all membership degrees representing the fuzzy set $A$. For any fuzzy subset, the relative frequency of the $j$-th fuzzy category of its non leaf node is shown in Eq. (7).

$$p_{ij}^{(k)} = \frac{M(T_i^{(k)} \cap T_i^{(n+1)} \cap X)}{M(T_i^{(n)} \cap X)} \tag{7}$$

In Eq. (7), $T_i^{(k)}$ represents a fuzzy subset; $X$ represents a non leaf node, and the fuzzy classification entropy of each fuzzy subset of $X$ is shown in Eq. (8).

$$FEntr_i^{(k)} = -\sum_{j=1}^{m} p_{ij}^{(k)} \log_2 p_{ij}^{(k)} \tag{8}$$

In Eq. (8), $1 \le k \le n$, $1 \le j \le m$, then the average fuzzy classification entropy of attributes on non leaf nodes is defined as shown in Eq. (9).

$$FEntr_k = -\sum_{j=1}^{m_k} \omega_i FEntr_i^{(k)} \tag{9}$$

In Eq. (9), $\omega_i$ represents the weight of the $i$-th attribute value. So the non assignability of classification on non leaf nodes is shown in the definition of Eq. (10).

$$Ambig_i^{(k)} = \sum_{j=1}^{m} (\pi_{ij}^{(k)} - \pi_{i,j+1}^{(k)}) \ln j \tag{10}$$

In Eq. (10), $m$ represents the number of fuzzy categories, so the average classification of non leaf nodes cannot be specified as shown in Eq. (11).

$$Ambig_k = -\sum_{j=1}^{m_k} \omega_i Ambig_i^{(k)} \tag{11}$$

When obtaining classification rules from uncertain information, fuzzy membership functions are often used to describe uncertainty. The fuzzy membership function characterizes the degree of uncertainty of the membership function through its parameters. This method can help handle data or situations that cannot be clearly classified into a certain category. The decision tree algorithm flow combining fuzzy logic is shown in Fig. 6.



Fig. 6. Decision tree algorithm flow combining fuzzy logic.

Fig. 6 shows the flowchart of the fuzzy decision tree algorithm. Compared with the standard decision tree algorithm, the steps for building the model have slightly changed, while the rest are roughly the same. When building a model, the dataset is first divided into different subsets on the ground of the selected features, with each subset corresponding to a sub node. It recursively constructs a subtree, recursively executing the above steps for each sub node until it reaches the leaf node or cannot be further divided. It designs a fuzzy set of nodes, and on the ground of the partitioning results and category labels, designs a fuzzy set representation of the current node, including fuzzy membership functions and membership degrees. This study evaluates the model on the ground of indicators such as accuracy, recall, F1 value, etc. The accuracy formula is shown in Eq. (12).

$$Accuracy = \frac{(TP+TN)}{(TP+TN+FP+FN)} \qquad (12)$$

In Eq. (12), $TP$ represents the true example; $TN$ represents a true negative example; $FP$ represents a false positive example; The recall rate formula is shown in Eq. (13).

$$Recall = \frac{TP}{TP+FN} \qquad (13)$$

In Eq. (13), $FP$ represents a false negative example; The recall rate represents the proportion of cases that the model correctly judges as positive, and can measure the model's ability to correctly recognize positive cases. The formula for F1 value is shown in Eq. (14).

$$F1 = \frac{2(Accuracy * Recall)}{(Accuracy + Recall)} \qquad (14)$$

## IV. ALGORITHM PERFORMANCE TESTING AND MODEL APPLICABILITY ANALYSIS

This chapter is separated into two sections. The first section mainly verifies the improvement effect of the decision tree algorithm and conducts comparative experiments with various similar algorithms. The second section mainly analyzes the applicability of the library intelligent acquisition decision-making model and applies it to actual library management.

### A. Performance Evaluation and Comparative Study of Fuzzy Decision Trees in Library Acquisition Decision-making

The experiment combines fuzzy sets to design a fuzzy decision tree on the ground of a decision tree. To study the algorithm performance and superiority of the model, Gradient Boosting Tree (GBT) and AdaBoos algorithm (DB) were introduced in the experiment to compare with the fuzzy decision tree proposed in the study. This experiment used PyTorch 1.8 software on the Windows 10 system platform to train three models 500 times using five different datasets, as shown in Fig. 7.

Fig. 7 shows the evaluation results statistics of three models in terms of recall and error rates. Fig. 7 (a) shows that the FDT model performs best in terms of recall rate, with various evaluation indicators reaching 96.45%, 92.64%, 87.92%, 91.17%, and 86.65%, respectively. In terms of false alarm rate, the FDT model also performs well, with significantly lower error rates than the other two models, which are 1.88%, 4.48%, 2.53%, 2.61%, and 2.86%, respectively; this means that the FDT model has a good effect in reducing false positives and can more accurately control the situation of erroneous reports. In summary, the FDT model performs well in terms of recall and false alarm rates, with high accuracy. The F1 value is an indicator that comprehensively considers the accuracy and completeness of the classifier, providing a more reliable performance evaluation under imbalanced datasets. The F1 values of the three models are shown in Fig. 8.



(a) Comparison of recall rates

(b) Comparison of false alarm rates

Fig. 7. Results of recall rate and false positive rate evaluation.

(a) F-measure of GBT



(b) F-measure of DB



(c) F-measure of FDT

Fig. 8.  F-measure for the three models.

Fig. 8 shows the results of training three models on five datasets. From the perspective of F1 values, the single GBT model is relatively low which is only 71.13% on dataset 1. Compared to this, the DB model exhibits good performance on most datasets, with a small and high average F1 difference, demonstrating its strong generalization ability. On the other hand, the F1 value of the FDT model is significantly higher than the other two models, which are 89.54%, 92.63%, 89.87%, 93.25%, and 94.79%, respectively. Overall, the average F1 value of the FDT model reaches 92.02%. The FDT model achieved the highest F1 value on all datasets, demonstrating its excellent performance in classification tasks. In addition, the AUC curves of each model were recorded in the experiment, as shown in Fig. 9.



Fig. 9.   AUC curves of three models.

Fig. 9 shows the statistical results of each model in terms of AUC values. In this figure, the horizontal axis represents the false positive rate, and the vertical axis represents the true positive rate. By observing Fig. 9, it can be observed that the ROC curve of the FDT model is always above the ROC curve of other models, indicating that its true positive rate performance is better under different false positive rates. In addition, on the ground of the enclosed area, i.e. AUC value, it can be concluded that the FDT model has the maximum AUC area of 0.681, while the AUC of the GBT and DB models are 0.517 and 0.463, respectively.

### B. Research on the Application of FDT Library Intelligent Acquisition Decision Model

The experiment first conducted extensive algorithm performance testing and comparative analysis on the FDT library intelligent acquisition decision-making model, fully proving the superiority of the model. However, to demonstrate the practical application value of the model, further research on its applicability is needed. By deploying the FDT model in an actual library environment and comparing it with traditional methods, the efficiency, accuracy, and application cost of the model were evaluated. Some of the library's procurement data are shown in Table I.

The experiment used the data in Table I to train the FDT library intelligent acquisition decision-making model and rule-based model. A rule-based model uses predefined rules and conditions to make interview decisions. For example, on the ground of factors such as user borrowing history and needs, establish a set of rules to determine whether to recommend a certain book to the user. And it uses False Positive Rate and False Negative Rate as evaluation indicators. False positive rate refers to the model mistakenly determining the proportion of resources that do not require interviews as those that require interviews; the false negative rate refers to the proportion of resources that the model mistakenly judges as not requiring interviews. The training results are shown in Fig. 10.

TABLE I.     PARTIAL PROCUREMENT DATA OF THE LIBRARY

| Registration no | Book attributes | Warehousing time | Classification number | Rice ($) |
|---|---|---|---|---|
| BC99851 | Storage | 2019.08.06 | B223.15 | 65.3 |
| BC98765 | Storage | 2019.08.06 | y | 56 |
| BC89732 | Circulate | 2019.08.06 | H319.4 | 58 |
| BC89875 | Storage | 2019.01.01 | H314 | 32.1 |
| BC89712 | Circulate | 2019.01.01 | H319.9 | 36.3 |
| BC88952 | Circulate | 2019.01.01 | I214.5 | 32.3 |
| BC84564 | Storage | 2019.09.01 | K512.4 | 50.2 |
| BC86832 | Circulate | 2019.01.01 | K512.4 | 12.3 |
| BC56465 | Circulate | 2019.07.01 | Y | 13.2 |



（a）The relationship between the number of training rounds and the variation of FP



（b）The relationship between the number of training rounds and the variation of FN

Fig. 10. The relationship between the number of training rounds and the variation of FP and FN.

Fig. 11. Rating of the recommendation effect of the FDT model by 48 students.

Fig. 10 shows the changes in FP and FN values of the FDT model and rule-based model as the number of training rounds increases. According to Fig. 10 (a), it can be observed that the false positive case (FP) value of the FDT model continues to decrease and tends to stabilize as the number of training rounds increases. On the contrary, the FP values of rule-based models are not very stable and show an increasing trend after 150 rounds of training. Further observation of Fig. 10 (b) shows that the error rejection rate (FN) of the FDT model rapidly decreases to about 0.1 as the training frequency gradually increases, and stabilizes at around 0.05 after 210 rounds of training. However, the error rejection rate of the rule-based library acquisition decision-making model ultimately stabilized at around 0.15 levels. These results emphasize the advantages and effectiveness of the FDT model in this classification problem. The experiment also randomly invited 48 middle school students from the library to rate the recommendation effect of the FDT library intelligent acquisition decision-making model. The relevant results are shown in Fig. 11.

Fig. 11 shows the rating of 48 students on the recommendation effectiveness of the FDT model and the rule-based library intelligent acquisition decision-making model. The figure shows that most students are satisfied with the recommendation performance of these two models, and the scores given are above 90 points. However, students rated the FDT library intelligent acquisition decision-making model higher than the rule-based model, at least 1.2 points higher. Specifically, the average score of the FDT model is 94.3, while the average score of the rule-based model is 91.8. In summary, most students hold a satisfactory attitude towards the recommendation effectiveness of FDT models and rule-based models.

## V. RESULTS AND DISCUSSION

In the research of intelligent procurement decision-making systems in libraries, traditional manual decision-making methods rely on experience and lack the ability to handle large datasets, making it difficult to quickly adapt to new needs and changes. Faced with this challenge, a series of improvements have been made to the standard decision tree, and a new fuzzy decision tree (FDT) model has been designed and implemented. This model aims to better handle uncertainty and ambiguity, and improve the accuracy and efficiency of decision-making. When evaluating the FDT model, several different datasets were used and compared with two popular algorithm models - Gradient Boosting Tree (GBT) and Adaptive Boosting Algorithm (DB) - for testing. The test results show that the GBT model has a relatively low F1 value on dataset 1, only 71.13%, while the DB model performs well on most datasets, but its average F1 value difference is not significant and lower than the FDT model. Although these two models perform well in certain aspects, they still have limitations in dealing with ambiguity and uncertainty. In contrast, the FDT model performs significantly better than the GBT and DB models on all test datasets. Specifically, the F1 values of the FDT model on five datasets were 89.54%, 92.63%, 89.87%, 93.25%, and 94.79%, respectively, with an average F1 value of 92.02%. This result highlights the significant ability of the FDT model to ensure high recall and accuracy, which is particularly crucial for libraries as it can reduce the risk of missing important books and avoid purchasing books that do not meet demand. The ROC curves of the FDT model always lie above the ROC curves of the GBT and DB models, indicating that the FDT model can maintain higher true positive rates at different levels of false positive rates. As an indicator of the overall performance of the model, the AUC value of the FDT model is 0.681, significantly higher than GBT's 0.517 and DB's 0.463. This result further demonstrates the superiority of the FDT model in distinguishing between different categories (books that need to be purchased and those that do not). In terms of practical application, 48 students evaluated the effectiveness of the FDT model, and the results showed that students generally rated the FDT library intelligent interview decision-making

model higher than rule-based models, with an average score of at least 1.2 points higher. The average score of the FDT model is 94.3 points, while the rule-based model is 91.8 points. This indicates that in practical applications, users are more satisfied with the recommendations provided by the FDT model, which may be because the FDT model can more accurately predict user needs and provide more personalized recommendations. However, despite the excellent performance of the FDT model in various aspects, research has also found a major drawback of this model: longer training time. This may be due to the model's need to evaluate and integrate a large number of fuzzy rules when processing data, resulting in an increase in computational complexity. The prolonged training process may limit the practicality of the model in application scenarios that require rapid updating of decision models to adapt to new situations. In the future, distributed computing resources can be utilized to allocate model training tasks to multiple computing nodes for parallel processing, significantly reducing training time. This requires the use of specialized distributed deep learning frameworks, such as the distributed version of TensorFlow, which is also an area that needs improvement in future research.

REFERENCES

[1] Fang B, Jiang M, Shen J. Deep Generative Inpainting with Comparative Sample Augmentation. Journal of Computational and Cognitive Engineering, 2022, 1(4): 174-180.

[2] Nimrah S, Saifullah S. Context-Free Word Importance Scores for Attacking Neural Networks. Journal of Computational and Cognitive Engineering, 2022, 1(4): 187-192.

[3] Guo Y, Mustafaoglu Z, & Koundal D. Spam Detection Using Bidirectional Transformers and Machine Learning Classifier Algorithms. Journal of Computational and Cognitive Engineering, 2022, 2(1), 5–9.

[4] Charbuty B, Abdulazeez A. Classification based on decision tree algorithm for machine learning. Journal of Applied Science and Technology Trends, 2021, 2(1): 20-28.

[5] Aldino A A, Sulistiani H. Decision Tree C4. 5 Algorithm for Tuition Aid Grant Program Classification (Case Study: Department of Information System, Universitas Teknokrat Indonesia). Jurnal Ilmiah Edutic: Pendidikan dan Informatika, 2020, 7(1): 40-50.

[6] Tangirala S. Evaluating the impact of GINI index and information gain on classification using decision tree classifier algorithm. International Journal of Advanced Computer Science and Applications, 2020, 11(2): 612-619.

[7] Li Q, Wen Z, He B. Practical federated gradient boosting decision trees//Proceedings of the AAAI conference on artificial intelligence.

[8] Nancy P, Muthurajkumar S, Ganapathy S, Santhosh Kumar, S. V. N., Selvi, M., & Arputharaj, K. Intrusion detection using dynamic feature selection and fuzzy temporal decision tree classification for wireless sensor networks. IET Communications, 2020, 14(5): 888-895.

[9] Ramya K, Teekaraman Y, Kumar K A R. Fuzzy-based energy management system with decision tree algorithm for power security system. International Journal of Computational Intelligence Systems, 2019, 12(2): 1173-1178.

[10] Elhazmi A, Al-Omari A, Sallam H, Mufti, H. N., Rabie, A. A., Alshahrani, M., ... & Arabi, Y. M. Machine learning decision tree algorithm role for predicting mortality in critically ill adult COVID-19 patients admitted to the ICU. Journal of infection and public health, 2022, 15(7): 826-834.

[11] Mariniello G, Pastore T, Menna C, Festa, P, & Asprone, D. Structural damage detection and localization using decision tree ensemble and vibration data. Computer-Aided Civil and Infrastructure Engineering, 2021, 36(9): 1129-1149.

[12] Li Q, Wu Z, Wen Z, & He, B. Privacy-preserving gradient boosting decision trees. Proceedings of the AAAI Conference on Artificial Intelligence. 2020, 34(1): 784-791.

[13] Sharma S, Obaid A J. Mathematical modelling, analysis and design of fuzzy logic controller for the control of ventilation systems using MATLAB fuzzy logic toolbox. Journal of Interdisciplinary Mathematics, 2020, 23(4): 843-849.

[14] Arji G, Ahmadi H, Nilashi M, Rashid, T. A., Ahmed, O. H., Aljojo, N., & Zainol, A. Fuzzy logic approach for infectious disease diagnosis: A methodical evaluation, literature and classification. Biocybernetics and biomedical engineering, 2019, 39(4): 937-955.

[15] Wang F, Song G. Bolt-looseness detection by a new percussion-based method using multifractal analysis and gradient boosting decision tree. Structural Health Monitoring, 2020, 19(6): 2023-2032.

[16] Sun L, Li Q, Fu S, & Li, P. Speech emotion recognition based on genetic algorithm–decision tree fusion of deep and acoustic features. ETRI Journal, 2022, 44(3): 462-475.

[17] Zhang J, Jia H, Zhang N. Alternate support vector machine decision trees for power systems rule extractions. IEEE Transactions on Power Systems, 2022, 38(1): 980-983.

[18] Prasenjit C. Model for selecting a route for the transport of hazardous materials using a fuzzy logic system. Vojnotehnički glasnik, 2021, 69(2): 355-390.

[19] Sahu R, Dash S R, Das S. Career selection of students using hybridized distance measure based on picture fuzzy set and rough set theory. Decision Making: Applications in Management and Engineering, 2021, 4(1): 104-126.

[20] Xiao F. A distance measure for intuitionistic fuzzy sets and its application to pattern classification problems. IEEE Transactions on Systems, Man, and Cybernetics: Systems, 2019, 51(6): 3980-3992.

# A Predictive Sales System Based on Deep Learning

Jean Paul Luyo Ballena, Cristhian Pool Ortiz Pallihuanca, Ernesto Adolfo Carrera Salas

Facultad de Ingeniería de la Universidad Peruana de Ciencias Aplicadas, Lima, Perú

*Abstract*—There are several techniques for predictive sales systems, in this study, a system based on different machine learning algorithms is developed for a trading company in Lima. As any company, it needs to be accurate in its sales calculations to manage the volume of production or product purchases. With the system, the trading company has a mechanism to order products from its supplier based on the predictions and estimates of the needs according to the projection of its sales. For the sales predictive system, Deep Learning technology and the neural network architectures GRU (Gated Recurrent Unit), LSTM (Long Short Term Memory) and RNN (Recurrent Neural Network) were used, 10 products were sampled, and the sales quantities of the last 12 months were obtained for the evaluation. The study found that the LSTM architecture excels in accuracy, significantly outperforming GRU and RNN in terms of Mean Absolute Percentage Error (MAPE), achieving an average MAPE of 7.07%, in contrast to the MAPE of 27.14% for GRU and the MAPE of 36.17% for RNN. These findings support the effectiveness and versatility of LSTM in time series prediction, demonstrating its usefulness in a variety of real-world applications.

*Keywords*—*Deep learning; neural network architectures; sales prediction; neural networks*

## I. INTRODUCTION

Effective sales management is essential to achieving business objectives but is often hampered by the availability of incomplete or outdated information, making it difficult to make informed decisions. Most of the business organizations heavily depend on a knowledge base and demand prediction of sales trends. The accuracy in sales forecast provides a big impact in business [1]. Despite the existence of conventional sales management systems, few take advantage of recent advances in artificial intelligence, and even fewer apply deep learning to improve accuracy in sales predictions. If the sales forecast is not accurate, situations of shortage or excess of stock may occur, which can have a direct and immediate impact on the profitability of the company. This effect is not only limited to the performance of profitability, but also to the customer service, therefore, service can be affected by an inefficient sales forecasting system. For example, if a customer is faced with a stock-out situation, they might decide to shop in a different company [2].

In this study, a Deep Learning based sales prediction system is proposed using three neural network architectures provided by the Brain.js library in JavaScript: GRU, LSTM and RNN. This research was conducted in a company located in Lima, using historical sales data to train and model each of these architectures. The research used 10 products specifically selected due to their relevance and analyzed data collected over a 12-month period to explore varied approaches to prediction.

To carry out this study, a series of key steps were followed. It started with the recording of historical sales data provided by the supplier. Then, different neural network architectures were used in JavaScript to perform sales analysis and prediction. Then, the data was loaded for neural network training, using the sales history as input data and the number of sales predicted for the next month as the expected output value. Subsequently, the neural network was modeled and trained with the architectures provided by the Brain.js library. Once the prediction for the next month was completed using each trained model, the corresponding predicted sales were displayed.

As part of the comparative analysis of the models, the predicted sales, the MAE (Mean Absolute Error) and MAPE (Mean Absolute Percentage Error) performance indicators and the inference time obtained from the three evaluated architectures are calculated. These steps allow to evaluate and select the most appropriate prediction model for the system, and to effectively apply the sales predictions.

Once the prediction for the next month was completed using each trained model, we show the corresponding predicted sales. As part of the comparative analysis of the models, we calculated the performance indicators MAE (Mean Absolute Error) and MAPE (Mean Absolute Percentage Error). Of the three architectures evaluated, the strongest result is the choice of LSTM due to its outstanding prediction accuracy, which significantly outperforms GRU and RNN in terms of MAPE, achieving a MAPE of 7.07%, in contrast to 27.14% for GRU and 36.17% for RNN.

This article is organized as follows: Section II shows the related works, where different neural network architectures are used for data analysis. It conceptualizes the different neural network architectures used in this study, as well as other relevant concepts used in the research. Section III explains the proposed methodology for the implementation of the Predictive Sales System, developing the following steps: Recording historical sales data, use of neural networks in JavaScript, Loading data for neural network training, Modeling and training the neural network, Model prediction and Comparison of results. In Section IV we show the results and discuss the findings, interpreting and analyzing their implications and limitations, and finally Section V concludes the paper.

In research [3], a novel and effective method for multi-channel retail demand forecasting is proposed. This method combines long-term memory (LSTM) and random forest (RF) neural networks to model complex temporal as well as regression relationships, which makes it more accurate than other forecasting methods. In addition, the proposed method has been shown to be statistically significantly better than other forecasting methods, including neural networks, multiple

regression, ARIMAX and LSTM. The proposed method also improves the interpretability of the model by ranking the relative importance of the explanatory variables, which makes it useful for decision making. An empirical evaluation was conducted on 192 time series and their corresponding information signals from two different channels: an online channel for 16 products and an offline channel covering the same 16 products sold in 11 different physical stores. For this reason, the analysis shows that the proposed method outperformed the other methods with a 95% confidence level in predicting the weekly demand for a product in a store. In conclusion, the proposed method is a valuable tool for multi-channel retailers in accurately and reliably predicting the demand for their products. This approach proves to be a valuable tool for retailers, providing accurate and reliable forecasts of the demand for their products, and contributing to our research on the performance of the LSTM algorithm.

In study [4], the authors used an LSTM neural network to evaluate its performance on point-of-sale data from a large retail chain. It was observed that bottom-up forecasts are more accurate than top-down forecasts when point-of-sale information is used for forecasting. The proposed framework helps to produce consistent short- and long-term forecasts for each level of decision making in a retail supply chain. This feature of the proposed framework helps align decision making at all levels of the organization and reduces the cost of decision misalignment. The authors conducted an empirical analysis using 141 different demand series for 10 different products offered through 1 online and 10 offline stores. Based on the results, it is safe to assume that at a confidence level of at least 95%, the bottom-up approach is better than the top-down approach in the case of retail demand forecasting. A top-down approach only generates a base forecast for the top principal node, and the top-down heuristic disaggregates the top-level forecast to obtain forecasts for all secondary nodes. In contrast, in the bottom-up approach, base forecasts are made for all lower-level nodes and then aggregated to obtain forecasts for all higher-level nodes. The article is included because of the similarity of the use of LSTM neural networks in the context of retail demand forecasting. This approach is critical to better understand how these networks can improve forecasting accuracy in these retail businesses.

In study [5], a meta-learning technique using deep convolutional neural networks, which have the feature of enhancing image recognition tasks and computer vision, is proposed for the purpose of automatically learning feature representation from sales time series. The meta-learner is used to automatically learn from a feature representation from raw sales time series data and then link the learned features to a common data set that is used to combine a set of methods from the base forecaster; the learned knowledge is used to select an optimal forecasting method for each time series according to its data features. Then, these characteristics are combined with a dataset to combine a set of forecasting methods to improve the accuracy of retail sales forecasting. The results indicate that the proposed meta-learner has a superior forecasting performance compared to several benchmark methods. The proposed meta-learner (M0) has significantly superior forecasting performance over all baseline forecasters and simple average combinations

of forecasters. On average, M0 has 3.2% improvements over the best performing base predictor. In addition, the meta-learner is particularly effective during promotion weeks, with a 5% improvement over the best-performing base predictor, compared to 1.6% improvements in non-promotion weeks. This approach helps us revolutionize retail sales forecasting. The deep convolutional neural network-based meta-learning technique proposed in this paper represents a valuable contribution to the field. Its ability to automatically learn time series characteristics and select optimal forecasting methods has demonstrated substantial improvements in forecasting accuracy. This not only benefits the trading company by enabling better planning and decision making, but also highlights the potential of applying deep learning and meta-learning in real-world business scenarios, as demonstrated in our sales forecasting research.

In study [6], the authors propose economic models based on exports and imports to forecast world trade using Deep Learning. The study presents a trade forecasting framework characterized by a theoretical integration of time series and economic structures, considering the uncertainty of international trade trends. The proposed approach demonstrates the power of trade data modeling using hybrid deep neural networks. In addition, the authors use prediction performance criteria to evaluate predictions using Root Mean Square Error (RMSE) and Mean Percentage Error (MAPE). The validity of this framework has been validated using commercial data from 10 countries: Brazil, Canada, China, France, Germany, India, Italy, Japan, UK, and USA. The predictive power of the proposed approach is compared with that of ten predictive models: ARIMA, ETS, TBATS, MLR, LASSO, RFR, XGB, SVR, MLP and HDL. The highest average values for (MAPE) and (RMSE), were observed in the time series model, reaching 2979 and 2933 for exports, and 4194 and 4475 for imports, respectively. In contrast, the models based on the economic structure of the time series showed lower average values for MAPE and RMSE, registering 2861 and 3198 for exports, and 3527 and 3716 for imports, respectively. The lowest MAPE values in their study stand at 2.861% for exports and 3.527% for imports, which contrasts significantly with our research, which obtains a minimum MAPE of 0.40%. These results underscore the ability of our model using LSTM to outperform in terms of accuracy the figures presented in the article, reinforcing the relevance and effectiveness of our methodology in accurately predicting future sales.

In study [7], the authors propose a study of electricity price estimation using deep learning approaches: an empirical study on Turkish markets in normal and Covid-19 periods. In recent decades, several methods have been used to estimate electricity bills. In the early days of these techniques, econometric models were used, but it was clear that these models were not successful enough to predict peak power prices. Furthermore, no transition period, such as COVID-19, is expected to affect electricity prices. However, machine learning and deep learning models are becoming increasingly important. Many researchers have shown that these models can also successfully predict extreme fluctuations in electricity prices. The proposed method is the original Transformer Encoder Decoder with Self-Attender (TEDSE) for electricity rate estimation. It is

organized in such a way that the model groups the data by a dummy variable and estimates the 24-hour electricity price of the Turkish electricity market. In this study, predictions were made using 12 different models based on Recurrent Neural Networks (RNN). The RNN-based estimates used to estimate electricity tariffs corroborated the results of the TEDSE model and underlined the success of the TEDSE model. However, due to market changes, these models individually could not accurately predict price estimates during COVID-19. For this reason, the TEDSE model forecasts sub-periods and allows electricity market participants to work more efficiently. This approach helps us understand the importance of deep learning, in particular, Recurrent Neural Networks (RNNs), in predicting complex phenomena such as electricity prices or as in our research on demand forecasting. In the study RNNs and advanced models, such as the Transformer Encoder-Decoder with Self-Attender (TEDSE), can provide more accurate forecasts even in challenging situations such as the COVID-19 period, where abrupt changes in markets are inevitable.

Table I presents a comparison of the works related to our research is made, where the different methods proposed by the authors are analyzed, the information on the data used and the time period of the study are detailed, in addition to the metrics evaluated to validate the models.

## II. CONTRIBUTION

The application of deep learning in sales prediction has been the topic of analysis on several research studies, as described in the related works in the previous section. The difference and contribution we make with our research is based on the evaluation of various architectures of recurrent neural networks, such as GRU, LSTM and RNN, which allows us to determine more accurately which architecture will provide us with the best metrics for predicting sales in a commercial company.

Based on the review of related works and the objective of our research, a broad conceptual framework is proposed to understand and define the fundamental concepts of Deep Learning, the architectures of neural networks and the company's need for an efficient sales management. In this context, it is important to clearly establish the key terms and concepts that will be discussed in this work, below we will find the essential definitions that will lay the foundation for solid understanding and deeper analysis in our research.

### A. Sales Management

Sales management in [8] refers to the comprehensive process of planning, coordinating, and controlling activities related to a company's sales force, with the objective of achieving optimal sales performance and reaching established business objectives. It involves developing effective strategies and tactics, as well as constantly monitoring and evaluating the performance of sales teams.

TABLE I. COMPARISON BETWEEN PROPOSED MODEL VS. RELATED WORKS

| Year | Author's References | Description of method / technique | Dataset Size | Performance Metric/s (value/s) |
|---|---|---|---|---|
| 2020 | Paper 1[3] | The proposed method is based on a state-of-the-art sequential deep learning method – long-short-term-memory networks (LSTM) and a machine learning method – random forest (RF). | Sales data of 16 products from online and offline stores over a period of 30 months. | The error measures employed were ARME, ARMAE, and ARMSE as proxies for the bias, accuracy, and variance of the evaluated forecasts. |
| 2020 | Paper 2[4] | Cross-Temporal Forecasting Framework (CTFF) to generate coherent forecasts at all levels of a retail supply chain, using a deep learning method, the long-short-term-memory network (LSTM). | Weekly sales data for 10 products (SKU) in a period of more than two years. | ARMAE, ARMSE, ARMAPE were used for the comparisons of forecasts errors across different approaches. |
| 2021 | Paper 3[5] | Meta-learning framework with automatic feature learning for retail product sales forecasting, based on convolution neural networks. | Sales data for 6 categories of products sold on 100 stores over a 153-week period. | Three error measures were used to compare the forecasting performance of the models: sMAPE, AvgRelMAE, MPE. |
| 2023 | Paper 4[6] | This study proposes a machine learning approach based on import and export economics statistics, the input features selection according to existing classical economic structure models under economics theories. | Data is collected from the of the export and import values of goods in USD records of 10 countries. All indices used in this study were adjusted using 2015 as the base year. | The methods used to evaluate the forecasts are the root mean square error (RMSE) and mean average percentage error (MAPE). |
| 2023 | Paper 5 [7] | The study aims to develop a prediction model for the electricity market of Turkey with more comprehensive Deep Learning models applying the TEDSE model. | In this study, 5-year of hourly historical data between 2017 and 2021 in the Turkish electricity market is used. | Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) statistics along with Mean Standard Error (MSE) were used to measure the model's performance. |
| 2023 | The model proposed | In this study, a Deep Learning based sales prediction system is proposed using three neural network architectures provided by the Brain.js library in JavaScript: GRU, LSTM and RNN. | The data for this study corresponds to historical sales records on 10 products in a period of 12 months. | To compare the performance of the models, MAE (Mean Absolute Error) and MAPE (Mean Absolute Percentage Error) are used. |

Source: Own elaboration

The importance of sales management lies in its ability to accelerate and improve the achievement of business objectives, as well as in the creation of competitive advantages and greater value in the marketplace. In addition, sales management facilitates the follow-up and monitoring of processes, identifying problems early and establishing training guidelines for the growth of the sales team. The three main objectives of sales management are to increase sales volume, contribute to

company profits and increase ROI, and achieve long-term sustainable growth. The system is dedicated to helping meet each of these objectives. Moreover, these objectives imply additional responsibilities for sales managers, who not only focus on sales, but must also manage and train a team of sales representatives to align them with organizational objectives. Effective sales management involves being informed about the market and making tactical decisions based on data and trend analysis.

In summary, sales management is a strategic process that seeks to maximize the performance of the sales force and achieve the company's business objectives by planning, coordinating, and controlling sales-related activities. Its importance lies in improving efficiency, generating competitive advantages, and ensuring the long-term growth of the organization.

### B. Deep Learning

Deep learning [9], is a subset of machine learning that relies on multi-layered neural networks to learn and make predictions from large volumes of data. Unlike conventional machine learning, deep learning can work with unstructured data, such as images or text, without requiring pre-processing to organize it into a specific format. Deep neural networks take advantage of algorithms such as gradient slope and back propagation to adjust their parameters and improve their accuracy over time. These algorithms automate feature extraction, which reduces the reliance on human intervention in manually defining relevant features. By eliminating the need for preprocessing and enabling learning directly from unstructured data, deep learning has driven the development of artificial intelligence applications in a variety of areas, such as speech recognition, fraud detection, autonomous driving and many more. Its ability to process complex information and learn hidden patterns has enabled significant advances in solving problems that were previously considered difficult or even impossible to tackle.

In short, deep learning is a branch of machine learning that uses multilayered neural networks to learn and make predictions from unstructured data. Its ability to automatically extract relevant features and its flexibility to work with various types of data have contributed to important advances in artificial intelligence.

### C. Neural Networks

In study [10], a neural network is a computational model inspired by the biological nervous system and consists of a collection of processing units called artificial neurons. These neurons are linked together by weighted connections that transmit and process information. Neural networks are used in many areas, including pattern recognition, computer vision, natural language processing and artificial intelligence in general. Their architecture consists of layers, including an input layer, one or more hidden layers and an output layer. Neural networks learn automatically through a training process in which connection weights are adjusted to improve their ability to make predictions and classifications. This process is based on learning algorithms such as supervised and unsupervised learning.

In short, a neural network is a computational model that imitates the nervous system, made up of artificial neurons connected to each other.

### D. JavaScript

In study [11], JavaScript is a programming language used to create interactive web pages. It enhances the user experience by enabling features such as social media updates, animations, and interactive maps. JavaScript was born as a technology intended for the browser, with the purpose of giving greater dynamism to web pages. Browsers were able to respond to user interactions and modify the layout of website content.

### E. RNN (Recurrent Neural Network)

The structure of an RNN in [12] consists of a series of repeating units connected in loops, also called "cells". Each repetitive unit receives input and produces output. There are also connections that can receive information from previous iterations. This allows information to flow recursively through the network and considers temporal dependencies.

### F. GRU (Gated Recurrent Unit)

GRU [13] is an efficient variant of SRRNs (standard recurrent neural network) used in natural language processing with its simple structure, it is easy to train and computationally efficient. It has proven its effectiveness in machine translation and speech recognition, being popular in machine learning and artificial intelligence.

### G. LSTM (Long Short-Term Memory)

In research [14], LSTMs are a special type of RNN whose main characteristic is that information can persist by introducing loops in the network diagram, which means that it can remember previous states and use this information to decide which will be the next one. LSTMs have a longer-term memory.

A great asset of LSTM [15], is its capacity to unravel intricate temporal structures and seize the fluid dynamics of systems that undergo time-driven fluctuations. By dissecting the inherent motifs and inclinations embedded within the data, LSTM architectures can reveal interconnections that may not be easily distinguished using established statistical or machine learning methodologies.

In comparison to LSTM, GRU [16] has a simplified cell structure that also operates based on a Gating system, but only has an update and reset gate. The main difference to LSTM is the circumstance that the cell state can be completely revised at each iteration and updated with short-term information via the reset gate. Instead, LSTM provides a mechanism that limits the change gradient that can be done at each iteration. Therefore, information is not completely gone with LSTM in contrast with GRU.

### H. MAE (Mean Absolute Error)

MAE [17] is a popular metric because the error value units match the predicted target value units. In MAE, different errors are not weighted more or less, but the scores increase linearly with the increase in errors. The MAE score is measured as the average of the absolute error values. The Absolute is a mathematical function that makes a number positive.

Therefore, the difference between an expected value and a predicted value can be positive or negative and will necessarily be positive when calculating the MAE.

### I. MAPE (Mean Absolute Percentage Error)

MAPE [18] is a relative error measure that uses absolute values to prevent positive and negative errors from canceling each other and uses relative errors that allows to compare forecast accuracy between time series methods.

Mean Absolute Percentage Error [19] is often used due to its very intuitive interpretation in terms of relative error. MAPE is relevant in finance, considering that profits and losses are often measured in relative values. It is also useful for calibrating product prices, customers are sometimes more sensitive to relative variations than to absolute ones. MAPE is frequently used when it is known that the quantity to be predicted remains well above zero. More generally, MAPE is well suited for forecasting applications, especially in situations where sufficient data is available.

### J. Inference Time

In research [20], inference time is the time it takes for a machine learning model to make predictions after training. It is important to optimize the model and hardware to meet the inference time requirements.

### III. METHOD

In this study, a sales management system was implemented using the Brain.js library in JavaScript. Fig. 1 presents the steps for the method used.

### A. Recording Historical Sales Data

The data for this study corresponds to historical sales records compiled from records provided by the company owner. We have focused on 10 products specifically selected because of their relevance. These products were accessed in the period between August 2022 and July 2023.

### B. Use of Neural Networks in JavaScript

To carry out sales analysis and prediction effectively, the Brain.js library was used in the JavaScript environment. This library offers several neural network architectures, such as GRU, LSTM and RNN, which were used to explore various approaches to prediction. This choice made it possible to investigate and select the most appropriate architecture to obtain accurate and relevant forecasting based on the patterns present in the historical data.

### C. Loading Data for Neural Network Training

At this stage, the training data for the neural network was prepared. An approach was followed in which the sales history was used as the input data, and the expected number of sales for the following month was used as the expected output value. In specific terms, the input consists of the 12 sales records that occurred in the past months (see Table II).

This structuring of the data allows the neural network to gain insight into the patterns and relationships underlying the

sales history and future projections. This approach is essential to achieve highly accurate predictions of future sales.

### D. Modeling and Training the Neural Network

At this stage, neural network training was carried out using the architectures provided by the Brain.js library: GRU, LSTM and RNN. Each architecture offers a different approach to modeling and learning relationships in the training data. Neural networks were modeled with the appropriate parameters and settings for each architecture, and then trained using the previously prepared data. During training, the neural networks were exposed to historical sales data to learn patterns and trends. As training progressed, the networks adjusted their weights and parameters to improve prediction accuracy. At the end of this stage, each of the three neural networks was ready to be used for prediction. Each architecture offers a unique perspective for making decisions related to future sales.

### E. Model Prediction

At this stage, we proceeded to use each of the trained models for the purpose of making predictions. The same procedure was followed for the three architectures available in Brain.js: GRU, LSTM and RNN. The prediction was carried out by supplying the sales value corresponding to each period in the dataset as input to the respective model. In this way, each architecture generated a sales prediction for the period in question. During this process, the inference time was recorded, reflecting the time required to complete each prediction, which made it possible to evaluate the performance of each model in terms of speed and efficiency. At the conclusion of each prediction, the sales estimate for each product was obtained and presented in the form of rounded integer values, which were displayed on the console along with the corresponding dates. This procedure was repeated for each of the three architectures, allowing a thorough comparison between the predictions generated by each model, as well as an evaluation of their accuracy and performance based on the training data.



Fig. 1. Flow of steps for the proposed system. Source: Own elaboration.

TABLE II.        MONTHLY SALES OF PRODUCTS (MONTHS 1-12)

| Products | Month | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| P-01 | 50 | 50 | 50 | 100 | 100 | 150 | 150 | 300 | 150 | 200 | 100 | 150 |
| P-02 | 50 | 50 | 100 | 100 | 100 | 150 | 100 | 300 | 150 | 150 | 100 | 100 |
| P-03 | 50 | 50 | 80 | 80 | 100 | 100 | 150 | 300 | 150 | 150 | 120 | 120 |
| P-04 | 200 | 250 | 250 | 300 | 250 | 300 | 300 | 200 | 150 | 100 | 100 | 100 |
| P-05 | 30 | 30 | 50 | 50 | 50 | 100 | 100 | 300 | 100 | 150 | 120 | 100 |
| P-06 | 30 | 20 | 20 | 50 | 50 | 50 | 80 | 200 | 100 | 100 | 100 | 150 |
| P-07 | 30 | 50 | 50 | 50 | 100 | 100 | 150 | 300 | 100 | 150 | 100 | 150 |
| P-08 | 50 | 50 | 40 | 40 | 80 | 100 | 100 | 300 | 100 | 150 | 150 | 100 |
| P-09 | 20 | 30 | 30 | 50 | 50 | 80 | 100 | 200 | 100 | 150 | 150 | 150 |
| P-10 | 50 | 50 | 100 | 50 | 50 | 100 | 100 | 300 | 100 | 150 | 100 | 200 |

Source: Own elaboration

### F. Comparison of Results

Once the prediction for the next month has been completed using each trained model, the corresponding predicted sale is displayed. In addition, as part of the comparative analysis of the models, MAE (Mean Absolute Error) and MAPE (Mean Absolute Percentage Error) performance indicators are calculated. These indicators are obtained by comparing the sales predicted by each model with the corresponding actual values of the data set, which provides a quantitative assessment of the accuracy of each model in its predictions and allows identifying and comparing its performance against the training data.

## IV. RESULTS AND DISCUSSION

This section provides the research result and discussion of the proposed models.

TABLE III.        RESULTS OF GRU (GATED RECURRENT UNIT)

| Products | GRU | | | |
|---|---|---|---|---|
| | Predicated Sale | MAE | MAPE | Inference Time (ms) |
| P-01 | 139 | 129 | 7,64% | 0,74 |
| P-02 | 61 | 41 | 38,53% | 0,80 |
| P-03 | 80 | 66 | 52,24% | 0,82 |
| P-04 | 88 | 38 | 12,30% | 0,90 |
| P-05 | 99 | 89 | 0,93% | 1,31 |
| P-06 | 50 | 20 | 75,17% | 1,06 |
| P-07 | 132 | 122 | 11,90% | 1,15 |
| P-08 | 99 | 89 | 0,54% | 1,06 |
| P-09 | 104 | 94 | 30,55% | 1,00 |
| P-10 | 117 | 97 | 41,57% | 0,86 |
| Average | | 79 | 27,14% | 0,97 |

Source: Own elaboration

To provide a clearer understanding of the results in Tables II, III and IV, data obtained from the three neural network architectures: GRU, LSTM and RNN, used to predict sales based on 10 products with a 12-month sales history as a dataset. The results focus on the sales forecast for the following month, analyzing the Mean Absolute Error (MAE), the Mean Absolute Percentage Error (MAPE) and the inference time for each product.

TABLE IV.        RESULTS OF LSTM (LONG SHORT-TERM MEMORY)

| Products | LSTM | | | |
|---|---|---|---|---|
| | Predicated Sale | MAE | MAPE | Inference Time (ms) |
| P-01 | 137 | 117 | 8,34% | 0,81 |
| P-02 | 101 | 81 | 0,60% | 0,85 |
| P-03 | 121 | 96 | 0,66% | 0,81 |
| P-04 | 83 | 33 | 17,01% | 0,72 |
| P-05 | 99 | 80 | 1,08% | 0,72 |
| P-06 | 136 | 96 | 9,21% | 0,82 |
| P-07 | 149 | 120 | 0,93% | 0,86 |
| P-08 | 100 | 90 | 0,40% | 0,75 |
| P-09 | 151 | 141 | 0,74% | 0,82 |
| P-10 | 180 | 90 | 31,70% | 0,88 |
| Average | | 94 | 7,07% | 0,80 |

Source: Own elaboration

TABLE V.        RESULTS OF RNN (RECURRENT NEURAL NETWORK)

| Products | RNN | | | |
|---|---|---|---|---|
| | Predicated Sale | MAE | MAPE | Inference Time (ms) |
| P-01 | 91 | 81 | 39,58% | 0,87 |
| P-02 | 93 | 73 | 7,48% | 0,78 |
| P-03 | 33 | 18 | 72,68% | 0,59 |
| P-04 | 115 | 65 | 15,30% | 0,80 |
| P-05 | 84 | 74 | 15,64% | 0,64 |
| P-06 | 77 | 67 | 48,73% | 0,60 |
| P-07 | 99 | 89 | 34,12% | 0,59 |
| P-08 | 66 | 56 | 34,42% | 0,61 |
| P-09 | 89 | 79 | 40,49% | 0,63 |
| P-10 | 94 | 74 | 53,27% | 0,60 |
| Average | | 68 | 36,17% | 0,67 |

Source: Own elaboration

Tables detailing the results for each architecture, GRU (see Table III), LSTM (see Table IV) and RNN (see Table V) provide clear and concise overview of the differences between them.

After analyzing the numerical results of the three evaluated architectures, the following results can be observed:

*1) Mean Absolute Percentage Error (MAPE):* The results in Table IV reveal that LSTM exhibits significantly lower MAPE values compared to GRU (see Table III) and RNN (see Table V). With an average MAPE of 7.07% for the 10 products, LSTM significantly outperforms GRU (27.14%) and RNN (36.17%). These findings highlight the remarkable precision of LSTM in predictions, demonstrating its ability to reduce the absolute error between the predicted mean value and the actual value.

Additionally, the incorporation of MAPE in this study provides a more comprehensive metric, as it not only considers the magnitude of the error but also the relative proportion to the actual value, offering a more detailed assessment of the predictive performance of each architecture.

*2) Mean Absolute Error (MAE):* Table V shows that the RNN architecture exhibits a lower average MAE (68) for the 10 products, surpassing GRU (see Table III) and LSTM (see Table IV), which have MAE values of 79 and 94, respectively. This suggests that, although RNN does not always predict with the same precision as other architectures, when it makes errors, the predictions do not differ significantly from the actual sales figures.

The inclusion of MAE in the analysis provides a valuable perspective on the average magnitude of errors, allowing for a more direct comparison between architectures in terms of absolute deviation of predictions from actual values.

*3) Inference Time:* The average inference times, presented in Tables III, IV and V are 0.97 milliseconds for GRU, 0.80 milliseconds for LSTM, and 0.67 milliseconds for RNN, respectively. While there are marginal differences between architectures in terms of inference speed, these results suggest that the choice of architecture should not be primarily based on inference time for this application. It is important to note that the effectiveness of each architecture should be evaluated considering both predictive accuracy and computational efficiency.

TABLE VI.    ARCHITECTURE PERFORMANCE COMPARISON

|  | GRU | LSTM | RNN |
|---|---|---|---|
| MAE | 79 | 94 | 68 |
| MAPE | 27,14% | 7,07% | 36,17% |
| Inference Time (ms) | 0,97 | 0,80 | 0,67 |

Source: Own elaboration

A key finding of the thorough analysis of the tables, is that LSTM stands out as the most favorable architecture for predicting product sales, as shown in Table VI. It not only achieves the lowest MAPE, indicating high prediction accuracy, but also maintains an acceptable MAE and a competitive average inference time, balancing speed with performance.

Our research has been characterized by working with three different architectures to enrich the literature and analyze different forms of data processing. This contributes to future research and obtaining a more efficient and timelier sales prediction, considering, that it is an important factor in companies knowing how much is going to be sold, which leads us to assess the quantities of items to be produced or sold, eliminating storage costs or overproduction, this being a success factor in any business.

LSTM has been successfully used in a variety of forecasting systems, ranging from retail sales prediction [3, 4, 5] to global trade forecasting in several major countries [6] and electricity price prediction in energy markets [7], highlighting its effectiveness and versatility in various forecast contexts. However, with technological advances in neural networks, it should have finer precision, noting that while the current minimum values of Mean Absolute Percentage Error and Mean Absolute Error are acceptable, they could be optimized.

## V.    CONCLUSION

The paper has analyzed several algorithms with the objective of improving the accuracy of future sales predictions, which supports the acquisition of products without the need to maintain a high warehouse stock and reduce unnecessary purchase costs per month. The results show that the LSTM architecture excels in terms of accuracy, implying a potential positive impact on the Lima trading company's operational efficiency and its ability to make informed strategic decisions.

Of the three architectures evaluated, the strongest result is LSTM, for its outstanding prediction accuracy that significantly outperforms GRU and RNN in terms of Mean Absolute Percentage Error (MAPE), achieving an average value of the 10 products at 7.07%, in contrast to 27.14% for GRU and 36.17% for RNN. In addition, an acceptable value of the Mean Absolute Error (MAE) with 94 and an average time of 0.80 milliseconds validated among the three architectures evaluated.

Finally, when considering prediction accuracy, Mean Absolute Percentage Error (MAPE) and inference time, the preferred choice is clearly LSTM. This is due to the importance of MAPE in this context. Although RNNs have advantages over MAE, LSTM remains the strongest choice for applications where prediction accuracy is crucial. LSTM focuses on measuring the relative accuracy of forecasts, which is essential for business decision making, effective comparisons, accurate inventory management, informed investments, efficient resource allocation and continuous improvement of business strategies. This underscores LSTM's ability to make highly accurate forecasts, which can be critical in applications where estimation accuracy is essential.

REFERENCES

[1] S. Cheriyan, S. Ibrahim, S. Mohanan and S. Treesa, "Intelligent Sales Prediction Using Machine Learning Techniques," International Conference on Computing, Electronics & Communications Engineering (iCCECE), Southend, UK, pp. 53-58, 2018, https://doi.org/10.1109/iCCECOME.2018.8659115.

[2] A.L.D. Loureiro, V.L. Miguéis, Lucas F.M. da Silva, "Exploring the use of deep neural networks for sales forecasting in fashion retail", Decision Support Systems, vol 114, pp. 81-93, 2018, https://doi.org/10.1016/j.dss.2018.08.010.

[3] Sushil Punia, Konstantinos Nikolopoulos, Surya Prakash Singh, Jitendra K. Madaan & Konstantia Litsiou, "Deep learning with long short-term memory networks and random forests for demand forecasting in multi-channel retail", International Journal of Production Research, vol. 58 no. 16, pp. 4964-4979, 2020, https://doi.org/10.1080/00207543.2020.1735666.

[4] Sushil Punia, Surya P. Singh, Jitendra K. Madaan, "A cross-temporal hierarchical framework and deep learning for supply chain forecasting". Computers & Industrial Engineering. vol. 149, 2020, https://doi.org/10.1016/j.cie.2020.106796.

[5] Shaohui Ma, Robert Fildes,- "Retail sales forecasting with meta-learning", European Journal of Operational Research, vol. 288, no. 1, 2021, pp. 111-128, https://doi.org/10.1016/j.ejor.2020.05.038.

[6] Cheng-Hong Yang, Cheng-Feng Lee, Po-Yin Chang, "Export- and import-based economic models for predicting global trade using deep learning", Expert Systems with Applications, vol 218, 2023, https://doi.org/10.1016/j.eswa.2023.119590.

[7] Mustafa Kaya, Mehmet Baha Karan, Erdinç Telatar, "Electricity price estimation using deep learning approaches: An empirical study on Turkish markets in normal and Covid-19 periods", Expert Systems with Applications, vol 224, 2023, https://doi.org/10.1016/j.eswa.2023.120026.

[8] Palacios D. " Sales management: what it is, why it matters and how to implement it". HubSpot, 2021, https://blog.hubspot.es/sales/gestion-de-ventas.

[9] "What is Deep Learning? | IBM", https://www.ibm.com/es-es/topics/deep-learning (accessed Jun. 25, 2023).

[10] "The neural network model | IBM Documentation", https://www.ibm.com/docs/es/spss-modeler/saas?topic=networks-neural-model (accessed Jun. 25, 2023).

[11] "What is JavaScript (JS)? | AWS", https://aws.amazon.com/es/what-is-javascript/#:~:text=AWS%20SDK%20para%20JavaScript%20es,Node (accessed Jun. 25, 2023).

[12] Zachary C. Lipton, John Berkowitz, Charles Elkan. "A Critical Review of Recurrent Neural Networks for Sequence Learning". arXiv preprint arXiv:1506.00019v4, 2015, https://doi.org/10.48550/arXiv.1506.00019.

[13] "Gated Recurrent Units – Understanding the Fundamentals | Data Science", https://datascience.eu/machine-learning/gated-recurrent-units-understanding-the-fundamentals/ (accessed Jul. 10, 2023).

[14] Muñoz González A., Salazar Guillen F. "LSTM neural networks application for the prediction of the daily movement direction problem for Bitcoin". In XI Congress of Applied Mathematics, Computational and Industrial - MACI 2021, Argentine Association of Applied Mathematics, 2021.

[15] Dongliang Li, Youyou Li and Zhigang ZHANG, "Automatic Extractive Summarization using GAN Boosted by DistilBERT Word Embedding and Transductive Learning" International Journal of Advanced Computer Science and Applications (IJACSA), vol. 14 no. 11, 2023. http://dx.doi.org/10.14569/IJACSA.2023.0141107.

[16] Benjamin Lindemann, Timo Müller, Hannes Vietz, Nasser Jazdi, Michael Weyrich, "A survey on long short-term memory networks for time series prediction", Procedia CIRP, vol. 99, 2021, pp. 650-655, https://doi.org/10.1016/j.procir.2021.03.088.

[17] Patrick Schneider, Fatos Xhafa, "Anomaly Detection and Complex Event Processing over IoT Data Streams, in Chapter 3 - Anomaly detection: Concepts and methods", 2022, pp. 49-66, https://doi.org/10.1016/B978-0-12-823818-9.00013-4.

[18] "MAPE Mean Absolute Percentage Error | Oracle® Fusion Cloud EPM Working with Planning", https://docs.oracle.com/cloud/help/es/pbcs_common/PFUSU/insights_metrics_MAPE.htm (accessed Jul. 10, 2023).

[19] Arnaud de Myttenaere, Boris Golden, Bénédicte Le Grand, Fabrice Rossi, "Mean Absolute Percentage Error for regression models", Neurocomputing, vol. 192, 2016, pp. 38-48, https://doi.org/10.1016/j.neucom.2015.12.114.

[20] Vergara J., "What are training and inference in artificial intelligence?" Neuroons, 2021, https://neuroons.com/es/que-son-entrenamiento-e-inferencia-en-inteligencia-artificial/ (accessed Jul. 10, 2023).

# Telemedicine and its Impact on the Preoperative Period

## A Systematic Review of the Literature

Raquel Elisa Apaza-Avila

Postgraduate Unit, Universidad Nacional Mayor de San Marcos, Lima, Perú

*Abstract*—The application of telemedicine has aroused a lot of interest in the field of chronic disease care, which is associated with clinical medicine. The aim of this research is to systematically evaluate the published evidence on telemedicine in the preoperative period. A systematic search was conducted over the last five years, excluding secondary research. Selection criteria were applied, obtaining 68 articles that met these criteria and quality criteria. The results show that the largest production is carried out in the United States and the United Kingdom, with collaboration between institutions and countries. The main use of telemedicine was in teleconsultation and telecounseling activities. In addition, the application of telemedicine in the preoperative period was made to a greater extent for general procedures without distinction of surgical specialty, oncological surgery and traumatology. An increased production observed can be related to the need for physical distancing due to the pandemic. Future research could include the co-occurrence of search terms, the impact of smartphones, NER terms, and the impact of polarity and objectivity on readers' choice of articles to read, share, and cite.

*Keywords—Telemedicine; digital health; e-health; preoperative care; preoperative period; systematic review*

## I. INTRODUCTION

The application of telemedicine has sparked interest in chronic disease care settings, which are generally encompassed in clinical disciplines. The application of telemedicine in surgical practice has become relevant in terms of the use of telesurgery. Reviews have been found about the use of applications in the prediction of mortality associated with surgery and the link with decision-making [1]. In this context, the question arises about the impact of telemedicine in the preoperative period.

There are a variety of experiences in the application of telemedicine in diverse health settings [2][3], which, despite the advantages regarding access to care, have also pointed out limitations such as difficulty in telephone access, omitted or erroneous information during data collection, and delay in the reporting of cases under investigation [4].

As Bokolo [5] points out, telemedicine and telehealth refer to the use of information and communication technologies embedded in software programs with high-speed telecommunications systems for the provision, management, and monitoring of health services.

In Peru, the definition of telemedicine is contained in the modifications of Telehealth Framework Law [6], as:

"The provision of remote health services in the components of promotion, prevention, diagnosis, treatment, recovery, rehabilitation and palliative care, provided by health personnel using ICTs, with the purpose of facilitating access to health services for the population".

Telemedicine could also be used to refer to the use of telecommunications for the remote provision of health services.

The World Health Organization (WHO) points out that telemedicine includes both diagnosis and treatment, as well as medical education, and that it is a technological resource that makes it possible to optimize health care services, saving time and money and facilitating access to distant areas for specialist care. In the context of the health crisis resulting from the pandemic, its use has become relevant [7].

Its applications include clinical practice and health education. Within clinical practice there are the following forms: Telediagnosis, Teleconsultation, Remote monitoring, medical meetings to obtain second opinions (Teleconference), Digital storage of data or medical records. Within the educational area, distance classes from medical centers (e-learning through videoconferencing) stand out [2] [4] [8].

Aspects related to the construction of bibliometric networks, polarity, objectivity, and subjectivity of the scientific production of telemedicine in the preoperative period have not been pointed out in the reviews on this topic.

In the current context, surgical activities have been suspended in most institutions around the world and considering the advantages of telemedicine in terms of timeliness of care [5], the aim is to investigate its usefulness in the field of preoperative activity.

This systematic review aims to explore the state of the art of telemedicine in the field of preoperative care.

In this vein, Section II covers the background and related works, where similar characteristics to the proposals of this work are specified. Section III is revision method which, details the methodology used in this document. Section IV delves into results and discussion which shows the compilation of studies and the data they generated, which are shown by graphs and tables to determine observations. To conclude, Section V, conclusions and future research, presents the recommendations reached because of the analysis of the

information obtained, as well as suggestions for scientific production on the subject addressed.

## II. BACKGROUND AND RELATED WORK

There are systematic reviews related to the application of telemedicine in the preoperative setting.

Research places telemedicine as a developing technology, and in the field of surgical practice, Sohn et al. [9] point to its use in plastic surgery and otolaryngology. In their review of telemedicine in the field of dermatological surgery, they point out that its application in preoperative consultation allows the planning of the intervention and increases access to care.

Bokolo et al., in its systematic review on the application of telemedicine and e-health technology in clinical services in response to the COVID-19 pandemic, points out the importance of the use of information and communication technologies (ICT) integrated with telecommunication software and systems for care, management and monitoring in patient care [5].

Asiri et al., on the other hand, in their review of the use of telemedicine in surgical care found that, for the most part, patients treated with this technology reported time savings and a reduction in the number of lost workdays as benefits [8].

However, not all reviews pointed to positive aspects. Moentmann et al., in their review on telemedicine in otolaryngology, noted that a negative aspect was the limitation of patient contact, although video-otoscopy is the most widely supported telemedical intervention limiting physical contact between otolaryngologists and their patients [9].

Kim et al. conducted a systematic review of research addressing the use of technology to intervene preoperatively on surgery anxiety in pediatric patients and their parents or guardians. They noted that the available literature is extremely heterogeneous and limits the ability to draw definitive conclusions about the effectiveness of technology-based interventions. In addition, the results showed that for this group of patients, tablets and manually operated devices with interactive capability may represent a viable option to address preoperative anxiety. However, they were unable to extrapolate these results to adults, with whom they had better results using videos [10].

More encouraging results are found with the reviews by Kolcun et al. [11] and Lu et al. [12]. In the first case, it highlights that the increase in the use of telemedicine has been favored by the crisis caused by the pandemic and represents an opportunity to continue developing this technology and validate its use in new fields. Their initial results show that this technology becomes a support for the interaction between doctors and patients during the need for social distancing, showing its usefulness for aspects that do not involve the need for physical contact. Aspects related to this point are indicated as, certain perioperative tasks (complementary patient education and postoperative surveys).

For Lu et al., in their review of the use of Short Message Service (SMS) and smartphones in surgical care, they conclude that applications of this type offer a sophisticated yet simple tool to improve perioperative health care, in addition to the need for a regulatory framework for communications [12].

Telemedicine is attracting attention in the healthcare sector, due to the diversity of interaction modalities that have been developed over the last decade, and which are becoming increasingly affordable for both patients and doctors. At this point, Shanbehzadeh et al. [13] highlight short message service, email and web portals, secure phone calls or VOIP, video calls, interactive mobile health applications (m-Health), remote patient monitoring, and video conferencing. At the same time, it points out that the synchronous modality through common social networks was the one that presented the highest percentage of use for clinical care. While data exchange activities using the store and forward service via secure messaging technology and pre-recorded media files had the least popularity.

## III. REVISION METHOD

The method used in this research is the systematic literature review (RSL), which is defined as a process of identifying, analysing and interpreting the existing scientific evidence on a topic, with the aim of providing answers to specific research questions.

The methodology used to develop the RSL in this paper is based on the document proposed by Kitchenham [14], who divides the whole process into three general parts: the planning of the review, the development, and the publication of results.

This research followed the phases defined by Kitchenham, as well as the activities that compose them. In the first phase, the research questions are specified, and the review protocol is developed, which is necessary to reduce the possibility of bias. In the second phase, the studies to be included in the research are identified, as well as the evaluation of their quality. Finally, in the third phase, the results obtained are detailed (see Fig. 1).



Fig. 1.  Development phases of the Systematic Literature Review (RSL). Translation of the systematic literature review process proposed by kitchenham.

## A. Problems and Objectives

When a systematic review of the literature is conducted, research questions are defined, which help in the extraction and analysis of data to meet the objectives of the research.

For this research, one general question and eight specific questions were posed.

The general question was:

What is the state of the art of Telemedicine and its impact on the Preoperative Period?

The objective of this study was to determine the current state of knowledge of the application of telemedicine in the preoperative period and to know the impact that this intervention generates in this period.

The specific questions and their objectives are shown in Table I.

TABLE I. CORRESPONDENCE BETWEEN RESEARCH QUESTIONS AND OBJECTIVES

| Question | Objective |
|---|---|
| RQ1. What are the most used and most relevant keywords by Number of Articles in telemedicine research and their impact on the preoperative period? | Determine which are the most used and most relevant keywords by Number of Articles in telemedicine research and in the preoperative period |
| RQ2: What is the relationship between the polarity of article titles and the frequency with which they are cited by other authors in telemedicine research and their impact on the preoperative period? | To determine the relationship between the polarity of article titles and the frequency with which they are cited by other authors in telemedicine research and its impact on the preoperative period |
| RQ3: What are the most productive institutions in the development of telemedicine and its impact on the preoperative period? | To determine which institutions are the most productive in the development of telemedicine and its impact on the preoperative period |
| RQ4. In which countries is telemedicine being applied most frequently in the preoperative period? | Determine where telemedicine is most commonly applied in the preoperative period |
| RQ5: Which means of publication are the main objectives for the production of research in the area of telemedicine in the preoperative period? | To determine the main means of publication for the production of research in the area of telemedicine in the preoperative period |
| RQ6. What are the types of telemedicine services that are most frequently provided in the preoperative period? | Determine which types of telemedicine services are most commonly provided in the preoperative period |
| RQ7. Which surgical specialties are most frequently applying telemedicine solutions in the preoperative period? | Determine which surgical specialties are most frequently applying telemedicine solutions in the preoperative period |
| RQ8. Which are the Articles whose Abstracts are characterized by their high Objectivity by year and country in research on telemedicine and its impact on the preoperative period? | To determine which articles whose abstracts are characterized by their high objectivity by year and country in telemedicine research and its impact on the preoperative period |

## B. Search Sources and Search Strategy

For this work, a bibliographic search was carried out using the most well-known search engines (see Table II).

TABLE II. SEARCH SOURCE

| Source |
|---|
| IEEE Xplore |
| Scopus |
| ARDI |
| ProQuest |
| ScienceDirect |
| ACM Digital Library |
| Wiley Online Library |
| Microsoft Academic |
| Springer |
| Google Scholar |

The table shows the search engines that were used to locate the research papers related to the topic of telemedicine and the preoperative period.

To determine the search terms, two well-known thesauri were used, the DeSC/MeSH for the terms related to telemedicine and the preoperative period, and the IEEE Xplore thesaurus also for the term telemedicine, and for the term methodology (see Table III).

TABLE III. IDENTIFICATION OF SEARCH TERMS

| Tesauro | Descriptor | Description |
|---|---|---|
| DeCS/MeSH IEEE Thesaurus | telemedicine digital health digital healthcare e-health m-health electronic health mobile Health | Independent Variable (A) |
| DeCS/MeSH | preoperative period preoperative care | Dependent Variable (B) |
| IEEE Thesaurus | methodology method model | Intervening Variable (C) |

The table shows the search engines that were used to locate the research papers related to the topic of telemedicine and the preoperative period.

The general equation was determined using the dependent, independent and intervening variables (see Fig. 2).

(telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health")

**AND**

("preoperative period" OR "preoperative care")

**AND**

(methodology OR method OR model)

Fig. 2. General search equation.

Equations based on the general equation were determined for each searcher (see Table IV).

## C. Identified Studies

The search yielded a total of 10,741 articles (see Fig. 3), to which filters related to the temporality of publication were applied, accessing articles from the last five years, as well as segregation by language, selecting those that were in English or Spanish. Subsequently, articles published in scientific journals and those peer-reviewed, as well as documents that were not duplicates, were selected.

TABLE IV.    EQUATIONS AND SEARCH SOURCES

| Search Source | Search Equation |
|---|---|
| IEEE Xplore | ("All Metadata":telemedicine OR "All Metadata":"digital health" OR "All Metadata":"digital healthcare" OR "All Metadata":e-health OR "All Metadata":m-health OR "All Metadata":"electronic health" OR "All Metadata":"mobile health") AND ("All Metadata":"preoperative period" OR "All Metadata":"preoperative care") AND ("All Metadata":methodology OR "All Metadata":method OR "All Metadata":model) |
| Scopus | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND ( methodology OR method OR model ) |
| ARDI | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND (method OR methodology OR model) |
| ProQuest | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND (methodology OR method OR model) |
| ScienceDirect | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND (methodology OR method OR model) |
| ACM Digital Library | [[All: telemedicine] OR [All: "digital health"] OR [All: "digital healthcare"] OR [All: e-health] OR [All: m-health] OR [All: "mobile health"] OR [All: "electronic health"]] AND [[All: "preoperative period"] OR [All: "preoperative care"]] AND [[All: method] OR [All: methodology] OR [All: model]] |
| Wiley Online Library | ""telemedicine" OR "digital+health" OR "digital healthcare" OR "e-health" OR "m-health" OR "electronic+health" OR "mobile health"" anywhere and ""preoperative period" OR "preoperative care"" anywhere and ""method" OR "methodology" OR "modeling"" anywhere |
| Microsoft Academic | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND (methodology OR method OR model) |
| Springer | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND (methodology OR method OR model) |
| Google Scholar | (telemedicine OR "digital health" OR "digital healthcare" OR e-health OR m-health OR "electronic health" OR "mobile health") AND ("preoperative period" OR "preoperative care") AND (methodology OR method OR model) |

The table shows the search engines that were used to locate the research papers related to the topic of telemedicine and the preoperative period.

## D. Exclusion Criteria

The following exclusion criteria were established for selecting articles:

CE1: Articles are more than five years old.

CE2: Articles are written in a language other than English or Spanish .

CE3: Articles followed peer review methodology and were not reported in a scientific journal.

CE4: The article did not propose a telemedicine solution or did not mention method or technique.

CE5: The article is not relevant to the objectives of the research.

SC 6: The article is not available, or the full text of the article is not available.

CE 7: The article is not unique.



Fig. 3.    Number of studies identified by search source.

## E. Selection of Studies

Initially, 10 741 articles were obtained, to which exclusion criteria were applied for the filtering and selection of the most relevant articles that provide better answers to the research questions posed (see Fig. 4 to Fig. 5).

As a result of this stage, a total of 68 articles were included (see Table V).

## F. Quality Assessment

To determine the final list of articles to be included in this research, criteria were applied to evaluate their quality.

Quality assessment criteria were determined for methodological characteristics and for substantive characteristics.

### 1) Methodological characteristics:

QA1: Are the objectives of the research clearly identified in the document?

QA2: Are research results clearly identified and reported?

### 2) Substantive features:

QA3: Does the research consider elective surgeries?

QA4: Is it possible to contact the principal investigator?

The full text for each document was analyzed and the criteria shown were applied to evaluate its quality and then conclude in maintaining the 68 articles.



Fig. 4.    PRISMA flowchart, on the application of criteria for the selection of articles.

TABLE V.        RESULT OF THE APPLICATION OF SELECTION CRITERIA

| Source | Initial Studies | Final Studies | % |
|---|---|---|---|
| IEEE Xplore | 3 | 1 | 1% |
| Scopus | 215 | 32 | 47% |
| ARDI | 421 | 4 | 6% |
| ProQuest | 16 | 1 | 1% |
| ScienceDirect | 1 828 | 3 | 4% |
| ACM Digital Library | 189 | 3 | 4% |
| Wiley Online Library | 128 | 6 | 9% |
| Microsoft Academic | 409 | 5 | 7% |
| Springer | 5 922 | 5 | 7% |
| Google Scholar | 1 610 | 8 | 12% |
| Total | 10 741 | 68 | 1% |

Note: Although ACM Digital Library provided the largest number of articles, the most relevant articles were obtained from Scopus.

### G. Data Extraction Strategies

At this stage, the final list of articles was used, from which the necessary information was extracted to answer the research questions (RQ1 to RQ8).

The data extracted from each article were: Article ID, Article Title, URL, Source, Year, Country, Number of Pages, Language, Type of Publication, Publication Name, Research Methodology, Author(s), Affiliation, Number of Citations, Abstract, Keywords, Conclusions/Discussions, Sample Size, RQ1, RQ2, RQ3, RQ4, RQ5, RQ6, RQ7, RQ8.

Not all articles answered all research questions.

The web and desktop application, Zotero, was used to manage data extraction (see Fig. 6).



Fig. 5.    The figure shows the result of the application of the search formula using the IEEE xplore, scopus, ARDI, ProQuest, sciencedirect, ACM digital library, wiley online library, microsoft academic, springer and google scholar search engines.



Fig. 6.    Document management with Zotero.

### H. Synthesis of Findings or Synthesis of Data

The information extracted for the Research Questions (RQ1 to RQ8) was tabulated and presented as quantitative data, using Excel, to statistically compare the various findings for each Research Question.

Certain patterns of research were found, as well as research directions that were carried out during the last few years.

Zotero was used for data management, while VOSViewer and Onodo were used for the analysis of bibliometric networks.

To determine objectivity and polarity, the Python program with the TextBlob library and the open access program CoreNLP v.4.3.2 were used.

## IV. RESULTS AND DISCUSSION

### A. Study Overview

Of the 68 articles included in the research, there has been a sustained increase in scientific production in the last two years (see Fig. 7).



Fig. 7. The figure shows the distribution of scientific production by year and source.

A variety of sources were searched, including those that are not common for health research publications. At this point, we compare the results with those obtained by Bokolo et al. [5], which we searched Google Scholar, PubMed, ScienceDirect, ProQuest, Springer, Sage, Taylor & Francis, IEEE Xplore, Wiley, ACM, Emerald, Inderscience, ISI Web of Science, and Scopus. The results are also compared with reviews of articles published in more well-known sources in the healthcare sector, such as the research conducted by Asiri et al., in its review on telemedicine in surgical care, in which MEDLINE, EMBASE, CINAHL and Science Direct were used to obtain articles [8].

Other reviews, such as that of Jonker et al., on e-health in the perioperative in older adults, included PubMed, EMBASE, CINAHL [77]. On the other hand, the team of Moentmann et al., in its review on telemedicine in otolaryngology, searched Embase, PubMed, and Web of Science, [9].

The number of articles included in these reviews is similar, except for the review by Jonker et al., in which the number of articles included was lower due to the delimitation of search criteria for the target group (older adults) [15].

As for the number of authors, they amounted to 436 in the 68 articles included. The number of authors varied in terms of the number of authors, with an average of six authors per publication. No collaborative relationships were found between the different research groups (see Table VI and Fig. 8).

A point to consider is related to the words that are most repeated in the titles (see Fig. 9). The most frequent words were identified as the words "preoperative", "study", "surgery", "telemedicine", "patients", "COVID-19" and "mobile", which are related to the search terms used.

TABLE VI. RESULT OF THE APPLICATION OF SELECTION CRITERIA

| N° | Source | Number of authors |
|---|---|---|
| 1 | IEEE Xplore | 8 |
| 2 | ProQuest | 7 |
| 3 | ScienceDirect | 7 |
| 4 | Scopus | 7 |
| 5 | ARDI | 6 |
| 6 | Google Scholar | 6 |
| 7 | Microsoft Academic | 6 |
| 8 | Springer | 6 |
| 9 | Wiley Online Library | 6 |
| 10 | ACM Digital Library | 4 |



Fig. 8. First authors and co-authors who formed research teams with a larger number of members.



Fig. 9. Word cloud of the titles of the articles included in the research.

### B. Answers to Research Questions

*1) RQ1.* What are the most used and relevant keywords by Number of Articles in telemedicine research and their impact on the preoperative period?

It is evident that the words that are most frequently used in medical articles are related to the search terms used. Table VII and Fig. 10 show the most frequently used keywords. It was consistent with other systematic reviews that include these keywords in their publications [5] [16], while other reviews, such as "surgical procedure", "satisfaction" y "monitoring" [8].

TABLE VII.    20 KEY WORDS MOST FREQUENTLY USED

| N° | Key Word | Number of Articles |
|----|----------|--------------------|
| 1 | telemedicine | 26 |
| 2 | humans | 19 |
| 3 | preoperative care | 17 |
| 4 | middle aged | 13 |
| 5 | female | 12 |
| 6 | male | 11 |
| 7 | aged | 9 |
| 8 | COVID-19 | 9 |
| 9 | prehabilitation | 7 |
| 10 | adult | 6 |
| 11 | telehealth | 6 |
| 12 | text messaging | 6 |
| 13 | patient satisfaction | 6 |
| 14 | perioperative care | 5 |
| 15 | surgery | 5 |
| 16 | ehealth | 5 |
| 17 | osteoarthritis | 5 |
| 18 | exercise | 4 |
| 19 | Mhealth | 4 |
| 20 | smartphone | 4 |



Fig. 10. Word Cloud of the keywords of the articles included in the research.

It has also been important to find co-occurrence between the keywords of the articles, such as "telemedicine", "COVID-19", "patient satisfaction", "preoperative care" and "prehabilitation", shown in Fig. 11. This result can provide

guidance on the impact that the pandemic has had on the development of telemedicine research in the preoperative period, which in turn is related in the articles to patient satisfaction and better preparation for surgery [5]. It should be noted at this point that the Named-entity recognition (NER) term search program in the titles of the articles, also identified the terms "COVID-19" and "COVID-19 Pandemic".



Fig. 11. Co-occurrance of keywords in the articles included in the study.

*2) RQ2:* What is the relationship between the polarity of article titles and the frequency with which they are cited by other authors in telemedicine research and their impact on the preoperative period?

Although no systematic reviews have been found that explore this point in the field of telemedicine, it is considered important to analyze the impact of this variable on readers.

As a result of the analysis of the titles of the articles, it was determined that, in general, titles with neutral polarity were the most cited, followed by those with positive polarity (see Fig. 12).



Fig. 12. Title polarity and citations.

On the other hand, the articles identified through the Microsoft and IEEE Xplore search engines showed greater neutrality in their writing (see Fig. 13).

Regarding the number of citations related to polarity and the search source in which the article was found, the highest frequency of citations is related to neutral titles extracted from Scopus (see Fig. 14).

Fig. 13. Title Polarity by Search Source



Fig. 14. The highest number of citations related to Scopus' neutral titles is observed.

*3) RQ3:* Which are the most productive institutions that establish collaborative networks in the development of telemedicine and its impact on the preoperative period?

Both public and private healthcare institutions, as well as those dedicated to research (universities, research groups) collaborated in the scientific production of telemedicine in the preoperative period. Six institutions produced two or more research articles (see Table VIII).

Collaboration between institutions is visualized in Fig. 15. Here we can see that the Technical University of Munich stands out. The articles that contributed the most to answering this question came from Scopus (see Table IX).

Some systematic reviews [1] [15] have pointed out the importance of collaboration between institutions and have included related experiences (first level and specialized centers, research institutions, universities, and hospitals) in their reviews.

*4) RQ4.* In which countries is telemedicine being applied most frequently in the preoperative period?

It is evident that publications related to telemedicine in the preoperative period have been carried out more frequently in the United States (49%) (see Table X and Fig. 16). This result is consistent with that described by M. Shanbehzadeh et al., in which the articles obtained were mostly (76.75%) carried out in this country [13].

Figures found in other reviews vary. Jonker et al. shows 28% [15], while Kolcun et al. 41.66% [11], with the United States occupying the first place in scientific production.

TABLE VIII. INSTITUTIONS THAT PUBLISH MOST FREQUENTLY ON TELEMEDICINE IN THE PREOPERTIVE PERIOD

| N° | Institution | Number of Articles |
|---|---|---|
| 1 | University of Michigan | 3 |
| 2 | The University of Melbourne | 3 |
| 3 | Dalhousie University | 2 |
| 4 | University of Cincinnati | 2 |
| 5 | Mayo Clinic College of Medicine | 2 |
| 6 | Vanderbilt University Medical Center | 2 |

TABLE IX. COLLABORATION NETWORKS BETWEEN INSTITUTIONS BY SEARCH SOURCE

| Source | Article | Quantity(%) |
|---|---|---|
| IEEE Xplore | [25] | 1 (2) |
| Scopus | [24] [28] [30] [33] [35] [37] [40] [42] [43] [44] [46] [48] [50] [60] [67] [72] [75] [78] [80] [81] [83] | 21 (51) |
| ARDI | [51] | 1 (2) |
| ProQuest | [19] | 1 (2) |
| ScienceDirect | [59] [73] | 2 (5) |
| ACM Digital Library | [49] [66] | 2 (5) |
| Wiley Online Library | [18] [71] | 2 (5) |
| Microsoft Academic | [63] [69] [77] | 3 (7) |
| Springer | [31] [53] [61] | 3 (7) |
| Google Scholar | [20] [34] [54] [74] [82] | 5 (12) |

Fig. 15. Collaborative networks between institutions that carry out research on telemedicine in the preoperative period.

TABLE X. SCIENTIFIC PRODUCTION BY YEAR AND COUNTRY

| Country | 2017 | 2018 | 2019 | 2020 | 2021 | Total (%) |
|---|---|---|---|---|---|---|
| Australia | | | 1 | 2 | 2 | 5(6) |
| Belgium | | | | 3 | | 3(4) |
| Canada | 2 | 1 | | 1 | 2 | 6(8) |
| China | | 1 | | 1 | 1 | 3(4) |
| Finland | | | | 1 | | 1(1) |
| France | | | | 1 | | 1 (1) |
| Germany | | | 1 | | | 1 (1) |
| India | | | 1 | | 1 | 2 (3) |
| Italy | | | | 1 | | 1 (1) |
| Mexico | | | 1 | | | 1 (1) |
| Netherlands | | | | 1 | 2 | 3 (4) |
| New Zealand | | 1 | | | | 1 (1) |
| Portugal | | | 1 | | | 1 (1) |
| Qatar | | | | | 1 | 1 (1) |
| Scotland | | | | 1 | | 1 (1) |
| Singapore | | | | 1 | | 1 (1) |
| Spain | | | 1 | | | 1 (1) |
| Sweden | | | 1 | | | 1 (1) |
| Taiwan | 1 | | | | | 1 (1) |
| United Kingdom | | 1 | 1 | 1 | 3 | 6 (8) |
| US | 6 | 5 | 7 | 12 | 9 | 39 (49) |

Regarding the establishment of collaboration networks with other countries, this research shows that the United States also leads this characteristic (see Fig. 17).



Fig. 16. Scientific production by country and year.



Fig. 17. Collaboration between countries on publications on telemedicine in the preoperative period.

*5) RQ5:* Which means of publication are the main objectives to produce research in telemedicine in the preoperative period?

Most of the publications correspond to journal-type articles (see Fig. 18).

This is consistent with other publications in the field of health [1] [5], in which the main input is publications of this type. It should be noted that some studies have only taken publications of this type as input, as in the research by Kolcun et al., which excludes publications such as "case reports", "technical reports" and "conference abstracts" [11].



Fig. 18. Research by type of publication.

*6) RQ6.* What are the types of telemedicine services that are available? Are they given more frequently in the preoperative period?

The main use of telemedicine in this period was in teleconsultation and telecounseling activities (see Fig. 19).

In this regard, the findings are consistent with the results of Asiri et al. [8], Kolcun et al. [11] and Shanbehzadeh et al., in which the majority use of telemedicine for teleconsultation and teleguidance activities was evidenced. Additionally, the use for telesurgery, tele-education and telemonitoring was reported [13].



Fig. 19. Types of service that are most frequently provided in the preoperative period.

As for the modality used, it was mainly characterized by being asynchronous (43%), however, it does not differ greatly from the synchronous modality (38%). 19% of publications use both modalities to provide telemedicine services (see Table XI and Fig. 20).

These results do not differ greatly from other reviews, in which both modalities were used [5], preferring videoconferencing for aspects related to diagnostic assessment [11].

TABLE XI.  TELEMEDICINE MODALITY USED

| Modality | Articles | Quantity (%) |
|---|---|---|
| Asynchronous | [20] [21] [25] [27] [28] [29] [30] [32] [36] [40] [48] [50] [51] [52] [54] [57] [58] [61] [62] [63] [66] [67] [69] [72] [73] [75] [77] [81] [83] | 29 (43) |
| Synchronous | [17] [18] [19] [22] [24] [26] [31] [34] [35] [37] [38] [43] [45] [47] [55] [56] [64] [68] [70] [71] [74] [76] [78] [79] [80] [82] | 26 (38) |
| Both | [23] [33] [39] [41] [42] [44] [46] [49] [53] [59] [60] [65] [84] | 13 (19) |

The communication channels used by the researchers varied, according to the activity carried out, but the use of videoconferencing and mobile applications stands out (see Table XII). These results coincide with studies carried out at the first level of care, such as the one conducted by A. C. Shah and S. M. Badawy in 2020 [85].



Fig. 20. Telemedicine modality most frequently used in the preoperative period.

TABLE XII.  COMMUNICATION CHANNELS IN TELEMEDICINE MOST FREQUENTLY USED IN THE PREOPERATIVE PERIOD

| Communication Channel | N° Articles |
|---|---|
| e-Form | 4 |
| e-Mail | 8 |
| Instant Messaging | 5 |
| Medical device | 2 |
| Mobile App | 16 |
| Phone Call | 15 |
| SMS | 7 |
| Smart device | 4 |
| Videoconference | 25 |
| Web | 9 |

No percentages have been placed in this table, since in about half of the publications they refer to the use of more than one communication channel at the same time.

These results could be linked to the emergence of new technologies associated with videoconferencing equipment and the expansion of smartphones [8].

Chen E.A. et al. [1], he clarifies this topic in his review "Smartphone applications in orthopedic surgery", mentions that the use of this equipment by physicians amounts to 90%, and performs a descriptive analysis of the use of mobile phones in the field of orthopedics, finding that their use in this field varied in capabilities from angular management to preoperative and gait quantification. And it concludes that as more advanced applications are developed, smartphones are likely to gain an increasing presence in both the operating room and clinical settings.

Something that we should also point out is that the articles included in this research point to interventions that used more than one communication channel (see Fig. 21).

*7) RQ7.* What are the surgical specialties that are most frequently applying telemedicine in the preoperative period?

Telemedicine investigations in the preoperative period were carried out to a greater extent without distinction of surgical specialty. The concentration of publications related to general preoperative management, oncological surgery, traumatology, general surgery, and neurosurgery is observed (see Fig. 22).

Fig. 21. Number of communication channels used in telemedicine activities during the preoperative period.

These results are in line with those published by Gachabayov et al., who addresses the issue of the role of telemedicine in surgical specialties during the pandemic and points out that most articles in the first six months were performed in orthopedic surgery followed by general surgery and neurosurgery, while in the second six months, urology and neurosurgery were the most productive, followed by transplantation and plastic surgery [86].



Fig. 22. The use of telemedicine in the preoperative period focuses on general procedures, oncological surgery, and traumatology.

*8) RQ8.* Which are the Articles whose Abstracts are characterized by their high Objectivity by year and country in telemedicine research and its impact on the preoperative period?

We observed that the production of articles during the first years was lower than during the last two years, but in recent years there has also been an increase in subjectivity in the abstracts of publications (see Fig. 23).



Fig. 23. Objectivity and Subjectivity of Abstracts

In terms of countries with highly objective summaries, the United States continues to lead (see Fig. 24).



Fig. 24. Abstracts with high objectivity by country.

## V. CONCLUSIONS AND FUTURE RESEARCH

This document has been an input and provided a statistical analysis on the application of Telemedicine in the Preoperative Period, through the extraction of data from a total of 68 articles published between 2017 and 2021. The highest percentage of identified studies was obtained from Springer, however, when applying the filtering and exclusion criteria, the highest percentage of included studies came from Scopus. It should be noted that the greatest use of telemedicine in this period is concentrated in teleconsultation and telecounseling services, as well as a greater scientific production with aspects related to general preoperative procedures, followed by those applied to oncological surgery and traumatology. There has also been an increase in production in recent years, probably due to the need for physical distancing due to the pandemic and the demand for activities in the surgical field.

For future research, it would be opportune to consider the co-occurrence of search terms, in this case, telemedicine with COVID-19 and preoperative care. It would also be a great contribution to analyze the impact smartphones have on preoperative care. Another relevant aspect would be to point out the use of NQER terms and the impact of polarity and objectivity on readers' choice of articles to read, share and cite.

## REFERENCES

[1] Chen, E. A., Ellahie, A. K., & Barsi, J. M. (2019). Smartphone applications in orthopaedic surgery: a review of the literature and application analysis: A review of the literature and application analysis.

Current Orthopaedic Practice, 30(3), 220–230. https://doi.org/10.1097/bco.0000000000000745

[2] Domingues, R. B., Mantese, C. E., Aquino, E. da S., Fantini, F. G. M. M., Prado, G. F. do, & Nitrini, R. (2020). Telemedicine in neurology: current evidence. Arquivos de Neuro-Psiquiatria, 78(12), 818–826. https://doi.org/10.1590/0004-282X20200131

[3] Pascual-de la Pisa, B., Palou-Lobato, M., Márquez Calzada, C., & García-Lozano, M. J. (2020). Efectividad de las intervenciones basadas en telemedicina sobre resultados en salud en pacientes con multimorbilidad en atención primaria: revisión sistemática. Atención primaria, 52(10), 759–769. https://doi.org/10.1016/j.aprim.2019.08.004

[4] Bertasso, C. P., Guerra, A. C. N., Pereira, F., Nakazato, L., Delatore, L. G., Anbar Neto, T., & Spadacio, C. (2021). Telemedicine in long-term elderly care facilities as "social accountability" in the context of Covid-19. Revista brasileira de educacao medica, 45(1). https://doi.org/10.1590/1981-5271v45.1-20200312.ing

[5] Bokolo, A. J. (2021). Application of telemedicine and eHealth technology for clinical services in response to COVID-19 pandemic. Health and Technology, 11(2), 1–8. https://doi.org/10.1007/s12553-020-00516-4

[6] Presidency of the Republic of Peru (2020). Legislative Decree N° 1490, Legislative Decree that strengthens the scope of telehealth. Official Gazette El Peruano of May 10, 2020. https://busquedas.elperuano.pe/dispositivo/NL/1866212-2

[7] M. Mihalj et al., "Telemedicine for preoperative assessment during a COVID-19 pandemic: Recommendations for clinical care", Best Practice & Research Clinical Anaesthesiology, vol. 34, núm. 2, pp. 345–351, jun. 2020, doi: https://doi.org/10.1016/j.bpa.2020.05.001.

[8] Asiri, S. AlBishi, W. AlMadani, A. ElMetwally, and M. Househ, "The Use of Telemedicine in Surgical Care: a Systematic Review," Acta Inform Med, vol. 26, no. 3, pp. 201–206, Oct. 2018, doi: 10.5455/aim.2018.26.201-206.

[9] M. R. Moentmann, R. J. Miller, M. T. Chung, and G. H. Yoo, "Using telemedicine to facilitate social distancing in otolaryngology: A systematic review," J Telemed Telecare, p. 1357633X20985391, Feb. 2021, doi: 10.1177/1357633X20985391

[10] Kim, J., Chiesa, N., Raazi, M., & Wright, K. D. (2019). A systematic review of technology-based preoperative preparation interventions for child and parent anxiety. Journal Canadien d'anesthesie [Canadian Journal of Anaesthesia], 66(8), 966–986. https://doi.org/10.1007/s12630-019-01387-8

[11] J. P. G. Kolcun, W. H. A. Ryu, and V. C. Traynelis, "Systematic review of telemedicine in spine surgery," Journal of Neurosurgery: Spine, vol. 34, no. 2, pp. 161–170, Oct. 2020, doi: 10.3171/2020.6.SPINE20863.

[12] K. Lu et al., "Use of Short Message Service and Smartphone Applications in the Management of Surgical Patients: A Systematic Review," Telemed J E Health, vol. 24, no. 6, pp. 406–414, Jun. 2018, doi: 10.1089/tmj.2017.0123.

[13] M. Shanbehzadeh, H. Kazemi-Arpanahi, S. Kalkhajeh, y G. Basati, "Systematic review on telemedicine platforms in lockdown periods: Lessons learned from the COVID-19 pandemic", J Edu Health Promot, vol. 10, jun. 2021, doi: 10.4103/jehp.jehp_1419_20.

[14] B. Kitchenham, "Procedures for Performing Systematic Reviews". [En línea]. Disponible en: https://bit.ly/3sZEd0H

[15] L. T. Jonker, M. E. Haveman, G. H. de Bock, B. L. van Leeuwen, and M. M. H. Lahr, "Feasibility of Perioperative eHealth Interventions for Older Surgical Patients: A Systematic Review," Journal of the American Medical Directors Association, vol. 21, no. 12, pp. 1844-1851.e2, Dec. 2020, doi: 10.1016/j.jamda.2020.05.035

[16] Sohn, G. K., Wong, D. J., & Yu, S. S. (2020). A review of the use of telemedicine in dermatologic surgery. Dermatologic Surgery, 46(4), 501–507. https://doi.org/10.1097/DSS.0000000000002230

[17] P. Hrishi, U. Prathapadas, R. Praveen, S. Vimala, and M. Sethuraman, "A Comparative Study to Evaluate the Efficacy of Virtual Versus Direct Airway Assessment in the Preoperative Period in Patients Presenting for Neurosurgery: A Quest for Safer Preoperative Practice in Neuroanesthesia in the Backdrop of the COVID-19 Pandemic!," J Neurosci Rural Pract, vol. 12, no. 04, pp. 718–725, Oct. 2021, doi: 10.1055/s-0041-1735824.

[18] K. F. Lee et al., "Mitigation of head and neck cancer service disruption during COVID-19 in Hong Kong through telehealth and multi-institutional collaboration," Head & Neck, vol. 42, no. 7, pp. 1454–1459, 2020, doi: 10.1002/hed.26226.

[19] Norcott et al., "Behaviours of older adults and caregivers preparing for elective surgery: a virtually conducted mixed-methods research protocol to improve surgical outcomes," BMJ Open, vol. 11, no. 10, 2021, doi: 10.1136/bmjopen-2020-048299.

[20] Naserian Mojadam, N. Nadeem, H. Beydoun, S. Abidi, A. Rizvi, and S. Abidi, "Preoperative Education System to Assist Patients Undergoing TAVI Surgery: A Digital Health Solution," Journal of Health & Medical Informatics, vol. 09, Jan. 2018, doi: 10.4172/2157-7420.1000313.

[21] H. Sadeghi et al., "Virtual reality and artificial intelligence for 3-dimensional planning of lung segmentectomies," JTCVS Techniques, vol. 7, pp. 309–321, Jun. 2021, doi: 10.1016/j.xjtc.2021.03.016.

[22] Joughin, S. Ibitoye, A. Crees, D. Shipway, and P. Braude, "Developing a virtual geriatric perioperative medicine clinic: a mixed methods healthcare improvement study.," Dec. 2021, doi: https://doi.org/10.1093/ageing/afab066.

[23] P. Norgan, M. L. Okeson, J. E. Juskewitch, K. K. Shah, and W. R. Sukov, "Implementation of a software application for presurgical case history review of frozen section pathology cases," in Journal of Pathology Informatics, Jan. 2017, vol. 8, pp. 3–3. doi: 10.4103/2153-3539.201112.

[24] Robinson, R. D. Sligth, A. K. Husband, and S. P. Slight, "Designing the optimal digital health intervention for patients' use before and after elective orthopedic surgery: ualitative study," Dec. 2021, doi: 10.2196/25885.

[25] S. Anusha et al., "Electrodermal Activity Based Pre-surgery Stress Detection Using a Wrist Wearable," IEEE Journal of Biomedical and Health Informatics, vol. 24, no. 1, pp. 92–100, Jan. 2020, doi: 10.1109/JBHI.2019.2893222.

[26] Rantala, M. Pikkarainen, and T. Pölkki, "Health specialists' views on the needs for developing a digital gaming solution for paediatric day surgery: A qualitative study," Journal of Clinical Nursing, vol. 29, no. 17–18, pp. 3541–3552, 2020, doi: 10.1111/jocn.15393.

[27] Shahrokni et al., "Electronic rapid fitness assessment: A novel tool for preoperative evaluation of the geriatric oncology patient," JNCCN Journal of the National Comprehensive Cancer Network, vol. 15, no. 2, pp. 172–179, 2017, doi: 10.6004/jnccn.2017.0018.

[28] Walter et al., "Improving the quality and acceptance of colonoscopy preparation by reinforced patient education with short message service: results from a randomized, multicenter study (PERICLES-II)," Gastrointestinal Endoscopy, vol. 89, no. 3, pp. 506-513.e4, 2019, doi: 10.1016/j.gie.2018.08.014.

[29] Huynh et al., "Patient and provider perceptions on utilizing a mobile technology platform to improve surgical outcomes in the perioperative setting," Journal of Surgical Oncology, vol. 123, no. 5, pp. 1353–1360, 2021, doi: 10.1002/jso.26406.

[30] M. Brennan et al., "Comparing clinical judgment with the MySurgeryRisk algorithm for preoperative risk assessment: A pilot usability study," Surgery (United States), vol. 165, no. 5, pp. 1035–1045, 2019, doi: 10.1016/j.surg.2019.01.002.

[31] A. Low et al., "A Real-Time Mobile Intervention to Reduce Sedentary Behavior Before and After Cancer Surgery: Usability and Feasibility Study," JMIR Perioperative Medicine, vol. 3, no. 1, p. e17292, Mar. 2020, doi: 10.2196/17292.

[32] Z. K. Christian et al., "Electronic Communication Patterns Could Reflect Preoperative Anxiety and Serve as an Early Complication Warning in Elective Spine Surgery Patients with Affective Disorders: A Retrospective Analysis of a Cohort of 1199 Elective Spine Patients," World Neurosurgery, vol. 141, pp. e888–e893, 2020, doi: 10.1016/j.wneu.2020.06.082.

[33] C. Conley et al., "The virtual pediatric perioperative home, experience at a major metropolitan safety net hospital," Dec. 2021, doi: https://doi.org/10.1111/pan.14179.

[34] A. H. S. Harris, A. C. Kuo, T. R. Bowe, L. Manfredi, N. F. Lalani, y N. J. Giori, "Can Machine Learning Methods Produce Accurate and Easy-to-Use Preoperative Prediction Models of One-Year Improvements in

Pain and Functioning After Knee Arthroplasty?", Journal of Arthroplasty, vol. 36, núm. 1, Art. núm. 1, 2021, doi: 10.1016/j.arth.2020.07.026.

[35] D. P. Lemanu et al., "Text messaging improves preoperative exercise in patients undergoing bariatric surgery," ANZ Journal of Surgery, vol. 88, no. 7–8, pp. 733–738, 2018, doi: 10.1111/ans.14418.

[36] D. S. Rubin et al., "Development and pilot study of an iOS smartphone application for perioperative functional capacity assessment," Anesthesia and Analgesia, vol. 131, no. 3, pp. 830–839, 2020, doi: 10.1213/ANE.0000000000004440.

[37] A. M. Delman et al., "Keeping the lights on: Telehealth, testing, and 6-month outcomes for orthotopic liver transplantation during the COVID-19 pandemic," Dec. 2021, doi: https://doi.org/10.1016/j.surg.2020.12.044.

[38] E. Layfield et al., "Telemedicine for head and neck ambulatory visits during COVID-19: Evaluating usability and patient satisfaction," Head & Neck, vol. 42, no. 7, pp. 1681–1689, 2020, doi: 10.1002/hed.26285.

[39] E. Piraux, G. Caty, G. Reychler, P. Forget, and Y. Deswysen, "Feasibility and preliminary effectiveness of a tele-prehabilitation program in esophagogastric cancer patients," Journal of Clinical Medicine, vol. 9, no. 7, pp. 1–14, 2020, doi: 10.3390/jcm9072176.

[40] F. De Mello-Sampayo, "Patients' out-of-pocket expenses analysis of presurgical teledermatology," Cost Effectiveness and Resource Allocation, vol. 17, no. 1, 2019, doi: 10.1186/s12962-019-0186-3.

[41] W. Fiona, R. Oloruntobi, L.-C. Roberto, and R. Tarannum, "The feasibility and effects of a telehealth-delivered home-based prehabilitation program for cancer patients during the pandemic," Dec. 2021, doi: https://doi.org/10.3390/curroncol28030207.

[42] "Implementation of telehealth is associated with improved timeliness to kidney transplant waitlist evaluation," Dec. 2021, doi: https://doi.org/10.1177/1357633X17715526.

[43] G. E. Halder et al., "A telehealth intervention to increase patient preparedness for surgery: a randomized trial," International Urogynecology Journal, 2021, doi: 10.1007/s00192-021-04831-w.

[44] G. Barugola, E. Bertocchi, and G. Ruffo, "Stay safe stay connected: surgical mobile app at the time of Covid-19 outbreak," Int J Colorectal Dis, vol. 35, no. 9, pp. 1781–1782, Sep. 2020, doi: 10.1007/s00384-020-03645-4.

[45] C. O. Hallet, F. J. Lois, D. O. Warner, J. A. Jastrowicz, J. L. Joris, and J. F. Brichant, "Short message service as a tool to improve perioperative follow-up of surgical outpatients: A before-after study," Anaesthesia Critical Care and Pain Medicine, vol. 39, no. 6, pp. 799–805, 2020, doi: 10.1016/j.accpm.2020.02.007.

[46] H. Al-Thani et al., "Implementation of vascular surgery teleconsultation during the COVID-19 pandemic: Insights from the outpatient vascular clinics in a tertiary care hospital in Qatar," PLOS ONE, vol. 16, no. 9, p. e0257458, 2021, doi: 10.1371/journal.pone.0257458.

[47] M. Herrera-Usagre et al., "Effect of a mobile app on preoperative patient preparation for major ambulatory surgery: Protocol for a randomized controlled trial," Dec. 2021, doi: 10.2196/10938.

[48] I. Idram, J.-Y. Lai, T. Essomba, and P.-Y. Lee, "Study on Repositioning of Comminuted Fractured Bones for Computer-Aided Preoperative Planning," in Proceedings of the 2017 4th International Conference on Biomedical and Bioinformatics Engineering, Nov. 2017, pp. 30–34. doi: 10.1145/3168776.3168801.

[49] J. L. Waterland et al., "Implementing a telehealth prehabilitation education session for patients preparing for major cancer surgery," Dec. 2021, doi: https://doi.org/10.1186/s12913-021-06437-w.

[50] J. Talevski et al., "Implementation of an electronic care pathway for hip fracture patients: a pilot before and after study," BMC Musculoskeletal Disorders, vol. 21, no. 1, pp. 1–7, Dec. 2020, doi: 10.1186/s12891-020-03834-w.

[51] J. Peralta et al., "Impact of a care delivery redesign initiative for vascular surgery," Journal of Vascular Surgery, vol. 71, no. 2, pp. 599-608.e1, Feb. 2020, doi: 10.1016/j.jvs.2019.03.053.

[52] K. Kulinski and N. A. Smith, "Surgical prehabilitation using mobile health coaching in patients with obesity: A pilot study," Anaesth Intensive Care, vol. 48, no. 5, pp. 373–380, Sep. 2020, doi: 10.1177/0310057X20947731.

[53] K. Reddy et al., "A Validation Pilot Study Comparing Telemedicine Images to a Face-to-Face Airway Exam for Conducting the Anesthesia Preoperative Airway Evaluation," Open Journal of Anesthesiology, vol. 11, no. 7, pp. 207–218, Jul. 2021, doi: 10.4236/ojanes.2021.117020.

[54] K. S. DeMartini et al., "Text Messaging to Reduce Alcohol Relapse in Prelisting Liver Transplant Candidates: A Pilot Feasibility Study", doi: https://doi.org/10.1111/acer.13603.

[55] K. Drummond, G. Lambert, B. Tahasildar, and F. Carli, "Successes and Challenges of Implementing Teleprehabilitation for Onco-Surgical Candidates and Patients' Experience: A Retrospective Pilot-Cohort Study," Nov. 2021, doi: https://doi.org/10.21203/rs.3.rs-1021190/v1.

[56] K. Taaffe, N. Zinouri, and A. G. Kamath, "Integrating simulation modeling and mobile technology to improve day-of-surgery patient care," in Proceedings of the 2016 Winter Simulation Conference, Arlington, Virginia, Dec. 2016, pp. 2111–2122. doi: 10.1109/WSC.2016.7822254.

[57] C. L. Kinman, K. V. Meriwether, C. M. Powell, D. T. G. Hobson, J. T. Gaskins, and S. L. Francis, "Use of an iPadTM application in preoperative counseling for pelvic reconstructive surgery: a randomized trial," Dec. 2021, doi: https://doi.org/10.1007/s00192-017-3513-2.

[58] L. Athlani, A. Chenel, R. Detammaecker, Y.-K. De Almeida, and G. Dautel, "Computer-assisted 3D preoperative planning of corrective osteotomy for extra-articular distal radius malunion: A 16-patient case series," Hand Surgery and Rehabilitation, vol. 39, no. 4, pp. 275–283, Sep. 2020, doi: 10.1016/j.hansur.2020.02.009.

[59] L. Newton and C. Sulman, "Use of text messaging to improve patient experience and communication with pediatric tonsillectomy patients," Dec. 2021, doi: https://doi.org/10.1016/j.ijporl.2018.07.048.

[60] M. Bendtsen, C. Linderoth, and P. Bendtsen, "Mobile Phone–Based Smoking-Cessation Intervention for Patients Undergoing Elective Surgery: Protocol for a Randomized Controlled Trial," JMIR Res Protoc, vol. 8, no. 3, p. e12511, Mar. 2019, doi: 10.2196/12511.

[61] S. Lee, A. Dana, and J. Newman, "Teledermatology as a Tool for Preoperative Consultation Before Mohs Micrographic Surgery Within the Veterans Health Administration," Dermatologic surgery : official publication for American Society for Dermatologic Surgery [et al.], vol. 46, no. 4, pp. 508–513, 2020, doi: 10.1097/DSS.0000000000002073.

[62] M. A. Woodward et al., "Tele-ophthalmic Approach for Detection of Corneal Diseases: Accuracy and Reliability," Cornea, vol. 36, no. 10, pp. 1159–1165, Oct. 2017, doi: 10.1097/ICO.0000000000001294.

[63] M. J. Heslin, J.-S. Liles, and P. Moctezuma-Velásquez, "The use of telemedicine in the preoperative management of pheochromocytoma saves resources," Dec. 2021, doi: 10.21037/mhealth.2019.08.04.

[64] M. T. Kemp et al., "Surgery Provider Perceptions on Telehealth Visits During the COVID-19 Pandemic: Room for Improvement," Dec. 2021, doi: https://doi.org/10.1016/j.jss.2020.11.034.

[65] M. W. Seward et al., "Weight loss before total joint arthroplasty using a remote dietitian and mobile app: study protocol for a multicenter randomized, controlled trial," Journal of Orthopaedic Surgery and Research, vol. 15, no. 1, pp. 1–8, Nov. 2020, doi: 10.1186/s13018-020-02059-w.

[66] M. A. Audette, T. Rashid, S. Ghosh, N. Patel, and S. Sultana, "Towards an anatomical modeling pipeline for simulation and accurate navigation for brain and spine surgery," in Proceedings of the Summer Simulation Multi-Conference, San Diego, CA, USA, Jul. 2017, pp. 1–12. doi: https://dl.acm.org/doi/10.5555/3140065.3140079.

[67] M. van der Velde et al., "Usability and preliminary effectiveness of a preoperative mHealth app for people undergoing major surgery: Pilot randomized controlled trial," JMIR mHealth and uHealth, vol. 9, no. 1, 2021, doi: 10.2196/23402.

[68] M. Mullen-Fortino, K. L. Rising, J. Duckworth, V. Gwynn, F. D. Sites, and J. E. Hollander, "Presurgical Assessment Using Telemedicine Technology: Impact on Efficiency, Effectiveness, and Patient Experience of Care," Telemedicine and e-Health, vol. 25, no. 2, pp. 137–142, 2019, doi: 10.1089/tmj.2017.0133.

[69] N. M. Vender, J. K. Plata, S. Rando, A. C. Draviam, D. Santa Mina, and M. Qadan, "Prehabilitation Telemedicine in Neoadjuvant Surgical Oncology Patients During the Novel COVID-19 Coronavirus Pandemic," Dec. 2021, doi: 10.1097/SLA.0000000000004002.

[70] N. V. Kamdar et al., "Development, Implementation, and Evaluation of a Telemedicine Preoperative Evaluation Initiative at a Major Academic Medical Center," Anesthesia and Analgesia, pp. 1647–1656, 2020, doi: 10.1213/ANE.0000000000005208.

[71] N. B. Seim et al., "Developing a synchronous otolaryngology telemedicine Clinic: Prospective study to assess fidelity and diagnostic concordance," in The Laryngoscope, 2018, vol. 128, pp. 1068–1074. doi: 10.1002/lary.26929.

[72] T. Osman et al., "PreAnaesThesia computerized health (PATCH) assessment: development and validation," BMC Anesthesiology, vol. 20, no. 1, 2020, doi: 10.1186/s12871-020-01202-8.

[73] P. W. Knapp, R. A. Keller, K. A. Mabee, R. Pillai, and N. B. Frisch, "Quantifying Patient Engagement in Total Joint Arthroplasty Using Digital Application-Based Technology," The Journal of Arthroplasty, vol. 36, no. 9, pp. 3108–3117, Sep. 2021, doi: 10.1016/j.arth.2021.04.022.

[74] Q.-L. Zhang, W.-P. Xie, Y.-Q. Lei, H. Cao, and Q. Chen, "Telemedicine usage via WeChat for children with congenital heart disease preoperatively during COVID-19 pandemic: a retrospective analysis," International Journal for Quality in Health Care, vol. 33, no. 2, p. mzab066, Jun. 2021, doi: 10.1093/intqhc/mzab066.

[75] P. N. Ramkumar et al., "Remote Patient Monitoring Using Mobile Health for Total Knee Arthroplasty: Validation of a Wearable and Machine Learning–Based Surveillance Platform," Journal of Arthroplasty, vol. 34, no. 10, pp. 2253–2259, 2019, doi: 10.1016/j.arth.2019.05.021.

[76] G. Rogers, "Using Telemedicine for Pediatric Preanesthesia Evaluation: A Pilot Project," Journal of Perianesthesia Nursing, vol. 35, no. 1, pp. 3–6, 2020, doi: 10.1016/j.jopan.2019.07.001.

[77] R. O. Alabi et al., "Novel use of telemedicine for corneal tissue evaluation in eye banking: Establishing a standardized approach for the remote evaluation of donor corneas for transplantation," Cornea, vol. 38, no. 4, pp. 509–514, Apr. 2019, doi: 10.1097/ICO.0000000000001848.

[78] S. C. Schallhorn, S. J. Hannan, D. Teenan, M. Pelouskova, and J. M. Schallhorn, "Informed consent in refractive surgery: In-person vs telemedicine approach," Clinical Ophthalmology, vol. 12, pp. 2459–2470, 2018, doi: 10.2147/OPTH.S183249.

[79] S. Fassas, E. Cummings, K. J. Sykes, A. M. Bur, Y. Shnayder, and K. Kakarala, "Telemedicine for head and neck cancer surveillance in the COVID-19 era: Promise and pitfalls," Head & Neck, vol. 43, no. 6, pp. 1872–1880, 2021, doi: 10.1002/hed.26659.

[80] S. Allsop, R. Fairhall, and J. Morphet, "The impact of pre-operative telephone support and education on symptoms of anxiety, depression, pain and quality of life post total knee replacement: An exploratory case study," International Journal of Orthopaedic and Trauma Nursing, vol. 34, pp. 21–27, 2019, doi: 10.1016/j.ijotn.2019.02.002.

[81] J. S. A. Song, L. Wozney, J. Chorney, S. L. Ishman, and P. Hong, "Design and validation of key text messages (Tonsil-Text-To-Me) to improve parent and child perioperative tonsillectomy experience: A modified Delphi study," International Journal of Pediatric Otorhinolaryngology, vol. 102, pp. 32–37, 2017, doi: 10.1016/j.ijporl.2017.08.029.

[82] W. Miller, S. Mohammadi, W. Watson, M. Crocker, and M. Westby, "The Hip Instructional Prehabilitation Program for Enhanced Recovery (HIPPER) as an eHealth Approach to Presurgical Hip Replacement Education: Protocol for a Randomized Controlled Trial," Dec. 2021, doi: 10.2196/29322.

[83] F. Zia et al., "Effects of a short message service (SMS) by cellular phone to improve compliance with fasting guidelines in patients undergoing elective surgery: a retrospective observational study," Dec. 2021, doi: 10.1186/s12913-020-06039-y.

[84] Q. Zuo, G. Zhang, and Y. Liu, "Health Education Using Telephone and WeChat in Treatment of Symptomatic Uterine Myoma with High-Intensity Focused Ultrasound," Dec. 2021, doi: 10.12659/MSMBR.911040.

[85] A. C. Shah and S. M. Badawy, "Telemedicine in Pediatrics: Systematic Review of Randomized Controlled Trials," JMIR Pediatrics and Parenting, vol. 4, no. 1, p. e22696, Feb. 2021, doi: 10.2196/22696.

[86] M. Gachabayov, L. A. Latifi, A. Parsikia, y R. Latifi, "The Role of Telemedicine in Surgical Specialties During the COVID-19 Pandemic: A Scoping Review", World J Surg, vol. 46, núm. 1, pp. 10–18, ene. 2022, doi: https://doi.org/10.1007/s00268-021-06348-1

# A Solution to Improve the Detection of the Nominal Value of the Financial Market: A Case Study of the Alphabet Stocks

Zhaohua Li[1]*, Xinyue Chang[2]

School of Management, Henan Institute of Technology, Xinxiang Henan, 453003, China[1]
School of Art and Design, Henan Institute of Technology, Xinxiang Henan, 453003, China[2]

*Abstract*—Given the regular occurrence of non-stationarity, non-linearity, and high levels of noise in time series data, predicting the value of stocks is a considerable difficulty. Traditional methods have the potential to enhance the precision of forecasting, although they concurrently introduce computational complexity, hence augmenting the probability of prediction inaccuracies. To effectively tackle a range of concerns, the existing body of research proposes a novel approach that combines a light gradient boosting machine, a machine learning methodology, with artificial bee colony optimization. In the context of the examined dynamic stock market, the proposed model demonstrated better efficiency and performance compared to alternative models. The recommended model exhibited optimal performance, characterized by a low error rate and high efficacy. The analysis utilized data about the stock of Alphabet over the period spanning from January 2, 2015, to June 29, 2023. The outcomes of the study provide evidence of the predictive accuracy of the proposed model in determining stock prices. The study's findings demonstrate how well the suggested model performs when it comes to correctly predicting stock prices. The proposed model presents a pragmatic methodology for evaluating and forecasting time series data about stock prices. The research's findings show that, in terms of forecast accuracy, the suggested model performs better than the methods currently in use.

*Keywords—Alphabet stock; machine learning; light gradient boosting machine; optimization; artificial bee colony algorithm*

## I. INTRODUCTION

Retail and institutional investors can buy and sell shares of publicly listed companies on stock exchanges to make money. It acts as a vital measure of a country's financial condition in general since it reflects firm performance and the business environment. The marketplaces where stocks are traded are the physical and online markets where both buyers and sellers come along. [1][2]. Investors, as well as traders, utilize a range of techniques to evaluate stocks and identify profitable opportunities. The investigation of stock price behavior has long piqued the curiosity of both academics and investors.

Consequently, several models have been created and tested to comprehend the underlying variables that affect the behavior of stock prices. Fama carried out one such study in 1965 [3]. This subject has been essential to developing the current knowledge of stock price behavior and is still an important area of study in the field of finance. The stock market is a complex system that fluctuates and is unpredictable. Given that pricing changes tend to be chaotic, noisy, nonparametric, nonlinear, nonstationary, and nonlinear, it is challenging for analysts to analyze and anticipate price changes accurately [4]. These characteristics imply that traditional statistical methods may not be sufficient for effective market analysis.

Therefore, a range of artificial intelligence and machine learning methods have been developed by researchers to circumvent these problems and improve the accuracy of stock market predictions. In contrast to traditional time series methods, reports based on machine learning can handle the stock market's complex, noisy, nonlinear, and unstable data to provide forecasts that tend to be more accurate [5]. Consequently, it has become the preferred technique for time series analysis in several areas [6][7]. Ensemble learning combines a variety of algorithms for machine learning to improve effectiveness, minimize mistakes, and increase accuracy [8]. This approach generates results that are more accurate and reliable than any one model could provide by integrating the results of several models that use various approaches and sets of features. The "boosting" machine learning method involves training many models in succession. Every new design aims to make the shortcomings of the prior one worse. Boosting is especially beneficial for classification and prediction issues since it may improve the precision of mediocre models. Due to their effectiveness, ensemble learning techniques like light gradient boosting machines (LGBM) and extreme gradient boosting (XGBoost) are often utilized. LGBM separates trees by the leaf and is quicker than XGBoost, which separates trees by depth [9][10]. Microsoft introduced LGBM, a machine-learning technique based on decision trees, by the end of 2017 [11]. In the machine learning field, especially in data science, LGBM has become quite well-liked because of its benefits of quick convergence time and less memory usage. It is commonly employed and discovered to be the successful answer in data mining contests [12]. Dhungana et al. [13] used LGBM to conduct an empirical investigation of S&P 500 price forecasting.

The utilization of artificial intelligence techniques in the stock market has gained significant traction due to its capacity to analyze vast quantities of data and discern patterns that are hard for humans to perceive. It is crucial to remember that these strategies' performance largely depends on how their initial parameter settings are created. Incorrect initialization will likely lead to forecasts and results that are untrustworthy [14][15]. Therefore, it is essential to thoroughly examine the

*Corresponding Author. Email: lzh104@126.com

initial configuration of the parameters when implementing artificial intelligence in stock trading. Consequently, a variety of optimization strategies have been developed to circumvent these restrictions, Like Aquila optimizer (AO) [16], Biogeography-based optimization (BBO) [17], Grey wolf optimization (GWO) [18], and Artificial bee colony (ABC) [19] and more can be used. The ABC optimizer is an optimization technique that exhibits excellent effectiveness, drawing inspiration from the behavioral patterns observed in honeybee colonies. The technique in question is classified under the swarm intelligence family and is extensively employed across diverse domains such as engineering, machine learning, and operations research. The aforementioned metaheuristic algorithm, which draws inspiration from nature, has a wide range of applications and possesses robust optimization capabilities [19].

This study tackles the formidable challenge of forecasting stock prices in time series data containing high levels of noise, non-stationary, and non-linearity. Driven by this undertaking, this research develops fundamental inquiries: In what ways can conventional forecasting techniques be refined to achieve higher levels of accuracy while preserving computational efficiency? Could the integration of an artificial bee colony optimization system with a light gradient enhancing machine present an innovative resolution to these obstacles? The primary aims of this research are to improve the accuracy of stock price predictions, reduce the complexity of computational tasks involved, and exhibit exceptional performance in this domain by utilizing the proposed fusion of methodologies. This work is significant due to the efficiency and efficacy of this model that has been demonstrated in dynamic stock markets. By employing Alphabet's stock data, the model's ideal performance, minimal error rate, and significant effectiveness, underscoring its superiority over current approaches have been demonstrated and its feasibility in practical scenarios.

The study contributed to forecasting a multivariate time series by outlining a brand-new training method based on metaheuristics for artificial bee colonies. This essay explores the field of ABC-LGBM hybrid stock price forecasting. The model was contrasted with several other models, including LGBM, GWO-LGBM, BBO-LGBM, and ABC-LGBM, to ascertain its correctness. The evaluation's findings offer insightful information on the effectiveness of the ABC-LGBM hybrid model and its potential as a stock price prediction tool. One of the analytical processes used in the inquiry was a thorough assessment of the data source and all of its pertinent features in the second part. The data was analyzed using several procedures, such as optimizer methods, evaluation metrics, and the LGBM model. The analysis results are provided and compared with results from other methodologies in the third part. The discussion part is provided in the fourth part. The investigation's results and recommendations are briefly described in the epilogue.

## II. MATERIALS AND METHODS

### A. Artificial Bee Colony Algorithm

The development of the Artificial Bee Colony (ABC) technique arose from an examination and emulation of the foraging habits of wild bees to address optimization problems through a mathematical model. The hired bees, observers, and scouts are the three different categories of honey bee ABC agents in the colony. Two groups of bees with an equal number of members each comprise the ABC algorithm's bee population. Onlooker bees are the other half, whereas employed bees comprise the first half. Within the context of the ABC algorithm, the determination of the bee's food supply location is considered an optimization problem, including several variables that require optimization. The assessment of the solution's fitness can be employed to characterize the excellence of the food supply about the objective function of this issue. In other words, finding the best answer is similar to the process that bees go through while searching for a suitable food source and the process of the ABC method shown in Fig. 1.

These are the specifics of the ABC algorithm: The initial solutions are produced at random and used as the locations of the bee agents' food sources. The bee agents are put through three primary cycles of repetition after startup. The process of selecting the most optimal and viable alternatives while actively avoiding subpar choices and continually revising and enhancing the workable solutions. The worker bees collectively select a novel potential food source location to enhance the efficacy of their solutions. They make their decision depending on the area surrounding the previously chosen food source. Utilizing Eq. (1), the location of the new food supply is determined.

$$v_{ij} = x_{ij} + \phi_{ij}(x_{ij} - x_{kj}), \tag{1}$$

where, $k \in \{1,2,3,..,SN\}$ and $k \neq i, j \in \{1,2,3,..,D\}$ are chosen at random indexes, $v_{ij}$ is a new feasible solution that is modified from its previous solution value $(x_{ij})$ based on a comparison with a randomly selected position from its neighboring solution $(x_{kj})$, and $\phi_{ij}$ is a random number between [-1,1] that is used to randomly adjust the previous solution to become an alternate solution in the next iteration. There is a positional difference between $x_{ij}$ and $x_{kj}$ in one specific dimension.

If a bee that is now employed encounters a new food source location with a higher fitness value, it will replace the previous food source position in its memory with the new one. Employed bees will share the nutritional advantages of their new food sources with the other bees when they return to their hive. The next step is for each observation bee to select one of the recommended food sources based on the fitness value calculated by the bees in use. The probability that a recommended food source will be picked is given in Eq. (2).

$$P_i = \frac{fit_i}{\sum_{i=1}^{SN} fit_i} \tag{2}$$

In which Fit $_i$ represents the fitness value of the food source $i$. $SN$ stands for the number of practical food sources.

The likelihood of an observer bee choosing a given food source is positively correlated with the fitness value of that food item. Once a food source has been chosen, the observer bees will proceed to the chosen food source and identify another potential food source location within the vicinity of the

earlier picked food source. The calculation and expression of the novel candidate food source can be determined using Eq. (1).

During the third stage, food-supplying positions lacking an enhanced fitness value will be discarded and substituted with a newly determined standing, assigned at random by a scout bee. This strategy aids in the prevention of suboptimal solutions.

The calculation for determining the initial random place selected by the scout bee is expressed by Eq. (3):

$$x_{ij} = x_j^{min} + \text{rand}[0,1]\left(x_j^{max} - x_j^{min}\right), \qquad (3)$$

The lower limit and higher limit of the food supply location in dimension j are represented by the $x_j^{min}$ and $x_j^{max}$, respectively.



Fig. 1. The diagram of the ABC optimizer.

The termination criteria, in this case, are based on the count of functional evaluation. The three major processes outlined above will be conducted iteratively until the predetermined number of function evaluations is reached. The cycle of bees can be shown in Fig. 2.

### B. Grey Wolf Optimization (GWO)

GWO method was initially created by Mirjalili et al. [20], which mimics the hunting habits and leadership hierarchy of the Grey Wolf (Canis lupus) in the wild. The magnificent grey wolves, who are members of the Canidae family, are extremely accomplished apex predators. They are powerful in their native environments due to their amazing hunting prowess and clever group strategies. These animals prefer to live in groups and

adhere to a strict social hierarchy. Depending on their characteristics, wolves may be divided into four different groups: alpha ($\boldsymbol{\alpha}$), beta ($\boldsymbol{\beta}$), delta ($\boldsymbol{\delta}$), and omega ($\boldsymbol{\omega}$). There is a common belief that within a pack of wolves, the alpha wolf is the member who possesses the most successful and efficient solution to any given problem or challenge that the pack may face. There is a common belief among some individuals that the alpha wolf, or the highest-ranking member in a wolf pack, possesses the most optimal solution to a given situation or problem and has authority over the whole pack of grey wolves. The beta and delta wolves come in second and third, respectively. The other wolves are from the omega group, which has a low rank. The alpha, beta, and delta wolves, who are the strongest, considerably help in hunting. The

responsibility of tracking, chasing, encircling, and attacking the target falls on the three wolves. Following are the three main steps of the grey wolf hunting process:

- The target is pursued and tracked down

- The prey is pursued, surrounded, and harassed until it stops moving

- The prey is aimed and attacked



Fig. 2. The illustration of the cycle of bees.

## C. Light Gradient Boosting Machine Algorithm

The construction of the model can be represented as follows since the LGBM approach is developed from decision trees. Given the data of training set $S = \{(x_i, y); f = 1, 2, \cdots, n; x_i \in \kappa^*, y \in R\}$ where $n$ is the number of samples containing features of $m$. To get an estimate, the forecasts generated by the decision trees are aggregated in the following manner:

$$y_i^{2a} = \sum_{p=1}^{p} f_p(x_0) \qquad (4)$$

With $f_p$ trees included as the trees, and there are $p$ total trees. To get $f_p$, the objective function below must be minimized.

$$f_p = \arg\min_{f} \sum_{i=1}^{n} L(y_i y^{\text{mop}}) + \Omega(f_p) \qquad (5)$$

The reduction function is marked as $L$, and the $\Omega$ parameter is for regularization. The specified value is provided by.

$$\hat{\Omega}(f_p) = a\,\mathrm{T} + \frac{1}{2}i\sum_{j=1}^{T} w_j^2 \qquad (6)$$

where, the penalty parameters for T leaves and the weight of the leaves, w, are $\alpha$ and $\lambda$, respectively. Assuming that, L represents the loss function and that it is a squared error. Then

$$L\left(y_i, \hat{y}_i^{LG(p-1)} + f_p(x)\right)$$

$$= \left(y_i - \hat{y}_i^{LG(p-1)} - f_p(x)\right)^2 = \left(r - f_p(x)\right)^{2^2} \qquad (7)$$

The residual, denoted as $r$, is utilized in the fitting process to generate the function $f_p$. Using a quadratic approximation function, the goal function at iteration $p$ is minimized.

$$f_P \simeq \arg\min_{f_p} \sum_{l=1}^{n} \left[ g f_P(x_i) + \frac{1}{2} h_i f_P^2(x_i) \right] + \Omega(f_P),$$

$$g_{\mathrm{I}} = \partial_j \omega_{(p-1)} L\left(y_l, \hat{y}_l^{LC(p-1)}\right), \qquad (8)$$

$$h_{\mathrm{I}} = \partial_{y^2\,_{(p-1)}^2} L\left(y_l, \hat{y}_l^{LC(p-1)}\right).$$

A new tree by minimizing the goal function, denoted as $f_p$, is produced. The decision tree partitions each node based on the criterion of the greatest data gain.

The variance gains for a node at position $s$ that divides feature $j$ are provided by:

$$Z_{j|O}(s) = \frac{1}{n_O}\left\{ \frac{\left(\Sigma_{\{x_i \in O: x_{ij} \leq s\}} g_i\right)^2}{n_{l|O}^j(s)} + \frac{\left(\Sigma_{\{x_i \in O: x_{ij} > s\}} g_i\right)^2}{n_{r|O}^j(s)} \right\}, \qquad (9)$$

## D. Data Preparation and Collection

To perform a comprehensive analysis, it is essential to include the trade volume as well as the Open, High, Low, and Close (OHLC) prices within a designated time frame. The dataset utilized in this research aims to enable the prediction of Alphabet stock market prices throughout an extensive temporal span, spanning from January 1, 2015, to the middle of 2023. Precise forecasting of stock prices has paramount importance for shareholders, financial professionals, and managers operating within the finance industry. This dataset contains the necessary historical stock price information and its related characteristics for conducting prediction analyses. The key sources of financial market information for the dataset are the stock exchanges and financial news outlets. Historical daily market share prices for Alphabet were collected for the chosen period. Between January 1, 2015, and mid-2023, this paper's dataset includes various variables about Alphabet stock shares for each trading day. The essential components encompassed within the context of stock market data are the specific date, the initial price at the commencement of the session of trading, the final price after the trading session, the maximum price attained by the shares throughout the day, the minimum price reached by the shares during the day, and the trading volume denoting the aggregate number of shares exchanged within the day. To ensure the quality and consistency of the data, rigorous data preparation procedures were implemented before conducting any predictive analysis. Data standardization was

performed to facilitate precise modeling and prediction. Normalizing data is the transformation of numerical parameters to a standardized range, often ranging from 0 to 1 or with a mean of 0 and a standard deviation of 1. When working on analytical or modeling projects, it's important to treat variables with different units or magnitudes equally. The use of diverse techniques, including scaling and normalization, is vital in the

process of data cleaning. These approaches play a crucial role in mitigating gradient mistakes and ensuring consistent training outcomes. Based on the information shown in Fig. 3, the dataset was divided into two parts: 80% of the data was allocated for training, while the remaining 20% was reserved for testing.



Fig. 3.   Dividing data for both training and testing.

### E. Assessment Criteria

The purpose is to evaluate the ability of combined models to make accurate predictions. This category comprises the root mean square error (RMSE), mean absolute percentage error (MAPE), mean squared error (MSE), and coefficient of determination ($R^2$).

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \qquad (10)$$

$$MAPE = \left(\frac{1}{n}\sum_{i=1}^{n}\left|\frac{y_i - \hat{y}_i}{y_i}\right|\right) \times 100 \qquad (11)$$

$$MSE = \frac{1}{N}\sum_{k=0}^{n}\binom{n}{k}(Fi - Yi)b^2 \qquad (12)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}} \qquad (13)$$

### III.   RESULTS

### A. Statistical Results

Table I presents an extensive compilation of statistical information about the dataset under consideration. The OHLC price and volume figures shown in this table provide a more clear depiction of the factual information. To conduct a more comprehensive and precise evaluation of the data, it is advisable to use statistical metrics such as the mean, amount, mean, std., min, median, max, and variance numbers.

### B. Results of the Models

Finding and evaluating the best hybrid algorithm for predicting stock prices is the main goal of this study.

Forecasting models have been developed, and complex variables that affect stock market trends have been investigated. The primary objective was to provide investors and analysts with valuable information to facilitate informed investment decision-making. Table II, Fig. 4 and Fig. 5 present a comprehensive assessment of the performance of each model, accompanied by an in-depth examination of its effectiveness.

A comprehensive data analysis assessment was conducted using four commonly employed metrics: RMSE, MSE, MAPE, $R^2$. These metrics are widely recognized for their ability to offer a precise assessment of the dependability, precision, and overall usefulness of the analysis. The efficacy evaluation of the LGBM model was conducted using the RMSE, MSE, MAPE, and $R^2$ metrics, comparing the model's performance with and without an optimizer. Enhanced understanding of the model's performance facilitates the ability to make well-supported judgments based on the obtained outcomes.

### C. Comparison with Other Works

In order to assess the accuracy of the ABC-LGBM model, the primary objective of this research was to examine its performance. To accomplish this, we conducted a comparative analysis of the $R^2$ values of our framework and other recent works. Upon meticulous examination, we determined that the ABC-LGBM exhibited superior performance to the other works, thereby substantiating the effectiveness of this framework. The specifics of this comparison are provided in Table III.

TABLE I.      DATASETS STATISTICAL SUMMARIES

|  | Open | High | Low | Volume | Close |
|---|---|---|---|---|---|
| amount | 2137 | 2137 | 2137 | 2137 | 2137 |
| mean | 70.05219 | 70.81457 | 69.3428 | 32.59751 | 70.09629 |
| Std. | 34.54605 | 34.97686 | 34.14654 | 15.6062 | 34.55914 |
| min | 24.66478 | 24.7309 | 24.31125 | 6.936 | 24.56007 |
| median | 58.4235 | 58.9 | 57.871 | 28.734 | 58.4095 |
| max | 151.8635 | 152.1 | 149.8875 | 223.298 | 150.709 |
| variance | 1193.43 | 1223.381 | 1165.986 | 243.5536 | 1194.334 |

## TRAIN



Fig. 4.    The values to evaluate each model during training.

# TEST



Fig. 5.    The values to evaluate each model during testing.

TABLE II.        THE RESULTS OF FORECASTING EVALUATION FOR THE BENCHMARKING METHODS

| MODEL / Metrics | TRAIN SET | | | | TEST SET | | | |
|---|---|---|---|---|---|---|---|---|
| | $R^2$ | *RMSE* | *MAPE* | *MSE* | $R^2$ | *RMSE* | *MAPE* | *MSE* |
| LGBM | 0.984 | 3.446 | 4.054 | 11.878 | 0.976 | 2.821 | 1.955 | 7.958 |
| GWO-LGBM | 0.988 | 2.940 | 3.286 | 8.646 | 0.982 | 2.449 | 1.709 | 5.997 |
| ABC-LGBM | 0.994 | 2.022 | 3.118 | 4.089 | 0.993 | 1.479 | 1.029 | 2.188 |

TABLE III.     COMPARISON OF THE MODEL IN COMPARISON TO THE OTHER RECENT WORKS

| References | Model | $R^2$ |
|---|---|---|
| [21] | LSTM | 0.977 |
| [22] | DNN LSTM | 0.972 |
| Present study | | 0.993 |

## IV.    DISCUSSION

This study aims to find and evaluate the best hybrid stock price prediction system. Complex stock market variables have been studied and forecasting algorithms constructed. Information for investors and analysts to make informed investment decisions was the main goal. Table II, Fig. 4 and Fig. 5 evaluate each model's performance and efficacy. A thorough data analysis was performed utilizing four popular

metrics: RMSE, MSE, MAPE, $R^2$. These criteria are extensively used to evaluate the analysis's dependability, precision, and usefulness. The LGBM model was evaluated for efficacy using RMSE, MSE, MAPE, and $R^2$ metrics, comparing performance with and without an optimizer.

Upon meticulous analysis of the training and test datasets, it was ascertained that the LGBM model, when implemented without the optimizer, yielded $R^2$ coefficients of 0.984 and 0.976 for the respective training and testing datasets. Furthermore, MAPE and MSE values obtained for the test dataset were notably lower at 1.955 and 7.958, respectively, in comparison to the corresponding values obtained during the training phase. Furthermore, RMSE values for the training and testing sets were 3.446 and 2.821, respectively. The integration of optimization techniques has significantly enhanced the efficacy of the LGBM model's performance. The usage of the GWO has yielded significant improvements, resulting in a rise in the $R^2$ value to 0.988 during the training phase and 0.982 during the testing phase. Furthermore, the MAPE and MSE values displayed a decrease in both the testing and training datasets.

Specifically, the MAPE values for the training and testing sets were 3.286 and 1.709, while the MSE values were 8.646 and 5.997, respectively. The RMSE values have decreased to 2.449 for the testing set and 2.940 for the training set. The

findings of this research demonstrate the efficacy of the optimization techniques employed in enhancing the efficacy of the LGBM model. The ABC-LGBM model demonstrated superior performance in comparison to the GWO-LGBM model. The values of $R^2$ obtained for the training phase and testing phases of the ABC-LGBM model were 0.994 and 0.993, respectively. Significantly, the training MAPE and MSE values exhibited a decrease to 3.118 and 4.089, respectively. Similarly, the testing MAPE and MSE values experienced a decrease to 1.029 and 2.188, respectively, thereby suggesting a significant enhancement in precision. Moreover, it can be observed that the RMSE values showed a decrease to 2.022 and 1.479 for the test and training datasets, accordingly. This reduction in error further substantiates the superior performance of the ABC-LGBM model in comparison to the GWO-LGBM model, as it demonstrates the ABC-LGBM's ability to provide accurate forecasts. The investigation demonstrates that the ABC-LGBM model exhibits superior effectiveness and efficiency compared to the GWO-LGBM model. The ABC-LGBM model exhibits a high level of efficacy, as seen in its outstanding $R^2$ scores of 0.994 for training and 0.993 for testing. The model had outstanding efficiency, as evidenced by its lowest MAPE and MSE testing values of 1.029 and 2.188, respectively. Findings suggest that the ABC-LGBM model exhibits a high level of precision and reliability.



Fig. 6. Forecasting curve training created with ABC-LGBM.

Fig. 7. Forecasting curve testing created with ABC-LGBM.

The ABC-LGBM model is widely regarded as a reliable and robust technique for generating highly precise market price forecasts. The stock share curves of Alphabet are presented in Table II, along with the corresponding figures, namely Fig. 6 and Fig. 7, which serve as evidence of the effectiveness of the model. In terms of precise prediction of stock prices, the ABC-LGBM model has superior performance compared to alternative models such as LGBM and GWO-LGBM. The ABC method is a highly effective technique that has been found to reduce price fluctuations significantly. By doing so, it has the added benefit of making trend prediction much easier and increasing overall model accuracy. This can be highly advantageous in a variety of industries where accurate forecasting is critical to success.

Additionally, the method is particularly useful in minimizing the impact of unforeseen events that might otherwise lead to significant market fluctuations. One of the distinctive features that sets the ABC-LGBM model apart from alternatives is its ability to learn from previous data effectively. In conclusion, the ABC-LGBM model has notable efficacy as a tool for predicting stock prices due to its high level of precision, accuracy, and capacity to assimilate information from previous datasets.

## V. CONCLUSION AND RECOMMENDATIONS

Forecasting stock prices is a complex endeavor characterized by a multitude of variables. The stock market is a complex and fluid process that is influenced by various factors, such as political events, societal dynamics, and economic conditions. To accurately assess the upcoming value of stocks,

it is imperative to consider a diverse array of financial statements, market trends, as well as other pivotal factors. Furthermore, the behavior of stocks can be significantly influenced by economic factors such as interest rates, inflation, and worldwide market conditions. Developing trustworthy models for forecasting is a challenging task due to the intricate nature and extensive array of components involved. The ability to generate precise forecasts necessitates a comprehensive comprehension of the capricious and nonlinear characteristics inherent in the marketplace. Fortunately, the ABC-LGBM model is a viable solution to these issues and has been proven to be dependable and precise. There are numerous ramifications of this research for the community. The proposed model begins by tackling a significant obstacle in the prediction of stock prices, a matter that holds considerable importance for financial analysts, investors, and individuals engaged in stock market operations. Through adeptly managing the intricacies of non-linearity, non-stationarity, and elevated noise levels present in time series data, the model furnishes a more dependable and precise instrument for predicting stock prices. Furthermore, the integration of artificial bee colony optimization with light gradient boosting machine integration represents a novel strategy that not only improves accuracy but also addresses the issue of computational intricacy. This has ramifications that extend beyond the realm of academia, as it provides pragmatic resolutions to the practical obstacles encountered by practitioners in the financial industry. The utilization of Alphabet's stock data for analysis enhances the study's practicality, given that Alphabet is a significant participant in the stock market. The research findings, which illustrate the

model's exceptional performance in comparison to alternative approaches, possess the capacity to impact the decision-making procedures of financial institutions and investors. In brief, this study not only enhances the theoretical comprehension of time series data analysis but also offers a practical and influential instrument for the field, enabling financial stakeholders to anticipate and navigate stock market trends more efficiently. The study examined the LGBM and GWO-LGBM models as part of its investigation into forecasting stock prices. However, the ABC optimizer technique, in combination with LGBM, produced the best outcomes. The study's dataset comprises Alphabet stock OHLC prices and volume from January 2, 2015, to June 29, 2023. Based on the findings of the inquiry, it has been determined that the ABC-LGBM model exhibits a high level of reliability and accuracy in predicting stock prices.

- Throughout the investigation, a comparative analysis was conducted to assess the accuracy and forecasting capacity of the ABC-LGBM model about other models. Based on the obtained data, it was consistently observed that the ABC-LGBM model exhibited superior performance compared to each of the other models. The $R^2$ value of 0.993 obtained after testing indicates a high level of accuracy, hence proving the precision of the predictions. The ABC-LGBM model consistently generated precise forecasts, as indicated by its minimal MAPE value of 1.029 and a low MSE value of 2.188. Overall, in terms of accuracy and efficacy, the ABC-LGBM model demonstrated outstanding results compared to the other approaches that were evaluated.

Recommendations:

Promotion of Hybrid Approach Adoption: Advocate for the implementation of the proposed hybrid approach, which integrates LGBM and ABC, on the grounds that it exhibits enhanced efficiency and performance in comparison to alternative models. Highlight the potential advantages that this model may offer in mitigating the difficulties presented by time series data's non-stationarity, non-linearity, and substantial levels of noise.

Application in Dynamic Stock Markets: It is recommended that the suggested model be implemented in dynamic stock markets as a means to improve the accuracy of predictions, particularly in situations where conventional approaches may introduce computational intricacy and possible errors. It is recommended to conduct additional validation and testing of the proposed model using a variety of stock datasets in order to evaluate its adaptability and efficacy across distinct market conditions. Investigate its functionality across a range of financial scenarios in order to ascertain its resilience.

Thorough Evaluation of Variables: Emphasize the significance of taking into account a wide range of financial statements, market trends, and critical factors in order to generate precise predictions regarding stock prices. Advocate for researchers and practitioners to adopt a comprehensive approach by incorporating a range of factors, such as economic conditions, political events, and societal dynamics.

Practical Significance: Elucidate the ways in which the ABC-LGBM model is effectively employed to tackle obstacles arising from the presence of noise, non-linearity, and non-stationarity within time series data. It is recommended to apply the model in practical financial contexts in order to assess its efficacy and influence on decision-making procedures.

Extension to Other Sectors: Advocate for an investigation into the feasibility of implementing the ABC-LGBM model in sectors other than finance, taking into account its capacity to offer practical solutions to computational complexities and forecasting obstacles. Promote collaborations across industries in order to capitalize on the capabilities of the model in various domains.

Sustained Comparative Analysis: Support the implementation of ongoing comparative analyses to evaluate the ABC-LGBM model's accuracy and predictive capability in comparison to emergent models and evolving methodologies. Motivate scholars to investigate its efficacy across diverse market environments and utilize an assortment of datasets.

REFERENCES

[1] Y.-H. Wang, C.-H. Yeh, H.-W. V. Young, K. Hu, and M.-T. Lo, "On the computational complexity of the empirical mode decomposition algorithm," Physica A: Statistical Mechanics and its Applications, vol. 400, pp. 159–167, 2014, doi: https://doi.org/10.1016/j.physa.2014.01.020.

[2] C. Zhang, J. Ding, J. Zhan, and D. Li, "Incomplete three-way multi-attribute group decision making based on adjustable multigranulation Pythagorean fuzzy probabilistic rough sets," International Journal of Approximate Reasoning, vol. 147, pp. 40–59, 2022, doi: https://doi.org/10.1016/j.ijar.2022.05.004.

[3] E. F. Fama, "Random walks in stock market prices," Financial analysts journal, vol. 51, no. 1, pp. 75–80, 1995.

[4] Y. S. Abu-Mostafa and A. F. Atiya, "Introduction to financial forecasting," Applied Intelligence, vol. 6, no. 3, pp. 205–213, 1996, doi: 10.1007/BF00126626.

[5] Y. Chen and Y. Hao, "A feature weighted support vector machine and K-nearest neighbor algorithm for stock market indices prediction," Expert Syst Appl, vol. 80, pp. 340–355, 2017, doi: https://doi.org/10.1016/j.eswa.2017.02.044.

[6] M. Zounemat-kermani, O. Kisi, and T. Rajaee, "Performance of radial basis and LM-feed forward artificial neural networks for predicting daily watershed runoff," Appl Soft Comput, vol. 13, no. 12, pp. 4633–4644, 2013, doi: https://doi.org/10.1016/j.asoc.2013.07.007.

[7] R. Bisoi, P. K. Dash, and A. K. Parida, "Hybrid Variational Mode Decomposition and evolutionary robust kernel extreme learning machine for stock price and movement prediction on daily basis," Appl Soft Comput, vol. 74, pp. 652–678, 2019, doi: https://doi.org/10.1016/j.asoc.2018.11.008.

[8] L. G. Kabari and U. C. Onwuka, "Comparison of bagging and voting ensemble machine learning algorithm as a classifier," International Journals of Advanced Research in Computer Science and Software Engineering, vol. 9, no. 3, pp. 19–23, 2019.

[9] T. R. Mahesh, V. Vinoth Kumar, V. Muthukumaran, H. K. Shashikala, B. Swapna, and S. Guluwadi, "Performance analysis of xgboost ensemble methods for survivability with the classification of breast cancer," J Sens, vol. 2022, pp. 1–8, 2022.

[10] Y. Zhou, W. Wang, K. Wang, and J. Song, "Application of LightGBM Algorithm in the Initial Design of a Library in the Cold Area of China Based on Comprehensive Performance," Buildings, vol. 12, no. 9, p. 1309, 2022.

[11] G. Ke et al., "Lightgbm: A highly efficient gradient boosting decision tree," Adv Neural Inf Process Syst, vol. 30, 2017.

[12] M. Ustuner and F. Balik Sanli, "Polarimetric target decompositions and light gradient boosting machine for crop classification: A comparative evaluation," ISPRS Int J Geoinf, vol. 8, no. 2, p. 97, 2019.

[13] S. Chakrabarty, P. Dhungana, and S. K. Sarada, "Application of Ensembles for Stock Index Price Prediction," 2022.

[14] Y. Han, P. Pan, H. Lv, and G. Dai, "A hybrid optimization algorithm for water volume adjustment problem in district heating systems," International Journal of Computational Intelligence Systems, vol. 15, no. 1, p. 39, 2022.

[15] M. Liang, S. Wu, X. Wang, and Q. Chen, "A stock time series forecasting approach incorporating candlestick patterns and sequence similarity," Expert Syst Appl, vol. 205, p. 117595, 2022.

[16] L. Abualigah, D. Yousri, M. Abd Elaziz, A. A. Ewees, M. A. A. Al-Qaness, and A. H. Gandomi, "Aquila optimizer: a novel meta-heuristic optimization algorithm," Comput Ind Eng, vol. 157, p. 107250, 2021.

[17] D. Simon, "Biogeography-based optimization," IEEE transactions on evolutionary computation, vol. 12, no. 6, pp. 702–713, 2008.

[18] S. Mirjalili, S. M. Mirjalili, and A. Lewis, "Grey wolf optimizer," Advances in engineering software, vol. 69, pp. 46–61, 2014.

[19] D. Karaboga, "Artificial bee colony algorithm," scholarpedia, vol. 5, no. 3, p. 6915, 2010.

[20] V. Chandran and P. Mohapatra, "Enhanced opposition-based grey wolf optimizer for global optimization and engineering design problems," Alexandria Engineering Journal, vol. 76, pp. 429–467, 2023, doi: https://doi.org/10.1016/j.aej.2023.06.048.

[21] Z. Jin, Y. Yang, and Y. Liu, "Stock closing price prediction based on sentiment analysis and LSTM," Neural Comput Appl, vol. 32, pp. 9713–9729, 2020.

[22] A. C. Nayak and A. Sharma, PRICAI 2019: Trends in Artificial Intelligence: 16th Pacific Rim International Conference on Artificial Intelligence, Cuvu, Yanuca Island, Fiji, August 26–30, 2019, Proceedings, Part II, vol. 11671. Springer Nature, 2019

# Analysis of the Financial Market via an Optimized Machine Learning Algorithm: A Case Study of the Nasdaq Index

Lei Wang, Mingzhu Xie*

School of Accounting and Finance, Anhui Xinhua University, Hefei Anhui, 230088, China

*Abstract*—The complex interaction among economic variables, market forces, and investor psychology presents a formidable obstacle to making accurate forecasts in the realm of finance. Moreover, the nonstationary, non-linear, and highly volatile nature of stock price time series data further compounds the difficulty of accurately predicting stock prices within the securities market. Traditional methods have the potential to enhance the precision of forecasting, although they concurrently introduce computational complexities that may lead to an increase in prediction mistakes. This paper presents a unique model that effectively handles several challenges by integrating the Moth Flame optimization technique with the random forest method. The hybrid model demonstrated superior efficacy and performance compared to other models in the present investigation. The model that was suggested exhibited a high level of efficacy, with little error and optimal performance. The study evaluated the efficacy of a suggested predictive model for forecasting stock prices by analyzing data from the Nasdaq index for the period spanning from January 1, 2015, to June 29, 2023. The results indicate that the proposed model is a reliable and effective approach for analyzing and forecasting the time series of stock prices. The experimental findings indicate that the proposed model exhibits superior performance in terms of predicting accuracy compared to other contemporary methodologies.

*Keywords—Stock market prediction; Nasdaq index; random forest; moth-flame optimization; MFO-RF*

## I. INTRODUCTION

Retail and institutional investors can purchase and profitably sell shares of publicly traded corporations on the stock market. The stock market is a vital gauge of a nation's overall economic health since it represents company success and the business climate. Exchanges, both physical and virtual, are places where buyers and sellers come together to exchange assets [1, 2]. Traders and investors use a variety of strategies to evaluate equities and identify profitable opportunities. Many models have been created and put to the test to comprehend the underlying elements that influence stock prices. Both investors and scholars have long been interested in the study of stock price behavior. Fama conducted a study on this topic in 1965, which is now a major area of study in finance and has had a big impact on the way how currently understand stock price behavior [3], analyzing the stock market is challenging due to its volatile and ever-changing character. Analysts find it difficult to effectively analyze and anticipate price changes due to the market's noisy, chaotic, dynamic, non-linear, non-

stationary, and nonparametric characteristics [3]. These qualities raise the possibility that conventional statistical techniques won't be enough for efficient stock market analysis.

To get around the challenges of making accurate stock market forecasts, academics have created several machine learning and artificial intelligence algorithms. These novel techniques aim to handle complex, noisy, chaotic, and non-linear stock market data more effectively than traditional time series methods, which often ignore the dynamic nature of the financial markets. Machine learning approaches employ sophisticated algorithms to analyze vast quantities of financial data and identify intricate patterns that may elude human observation. These methodologies have the capacity to provide more precise predictions through the analysis of an extensive array of data sources, including news articles, social media postings, and financial information. Furthermore, machine learning algorithms have the capability to consistently acquire knowledge and adapt to novel data, hence improving their predictive abilities as time progresses. In general, the utilization of artificial intelligence and machine learning has the potential to significantly enhance the precision of stock market prediction, providing investors with valuable information on market fluctuations and facilitating prudent investment choices [4].

Decision trees (DT) are a widely utilized machine learning approach that finds frequent application in both classification and regression tasks. This basic and comprehensible method is employed to construct a tree-like model that represents various choices and their corresponding probable outcomes [5]. The method partitions the data into subsets based on the values of the input features iteratively until a specified stopping condition is satisfied. Every partition is a node in the tree, and every node contains a decision rule that determines which feature will be separated first. DTs provide several benefits in comparison to other machine learning approaches. First of all, because they can handle both continuous and categorical data, they are suitable for a wide range of applications. Secondly, they are easy to use and need little feature engineering or data preparation.

Moreover, they are easily interpreted, which simplifies the process of understanding how the model arrived at a certain prediction. Despite its advantages, DTs have some disadvantages. For example, the tree may be susceptible to overfitting if it is deep and complex. Additionally, owing to their sensitivity to even minute changes in data, they could

generate different trees for different training sets. To get around these problems, many DT variants have been developed; one of the best is random forest (RF).

RF employs an ensemble learning methodology to combine many decision trees, resulting in a robust and accurate model [6]. RF is a versatile machine-learning algorithm that can be applied to both classification and regression tasks. It has exceptional performance in handling large-scale datasets with high-dimensional features. The technique operates by constructing a collection of decision trees, whereby each tree is trained on a randomly selected portion of the characteristics and data. During the training process, every tree within the forest generates a forecast, and the collective projections of all the trees culminate in the ultimate prediction. This methodology enhances the model's robustness against noise and outliers, hence mitigating the risk of overfitting. When juxtaposed with alternative machine learning methodologies, RF offers several advantages. The versatility of this method is shown in its ability to effectively handle both continuous and categorical data, rendering it highly suitable for a wide range of applications.

Furthermore, this approach does not need an extensive feature engineering or data preprocessing, making it very straightforward to implement. In conclusion, the system exhibits a high degree of scalability and possesses the ability to effectively process extensive datasets comprising millions of samples and numerous attributes [6]. In their study, Park et al. [7] developed a comprehensive framework for predicting stock market trends by combining long short-term memory and random forest techniques. To evaluate the effectiveness of their proposed approach, the researchers utilized three prominent global stock indexes and incorporated 43 financial, technical indicators. In their study, Basher et al. [8] employed the RF algorithm to predict Bitcoin prices. Their findings indicate that the RF algorithm outperforms logit models in accurately anticipating trends in both Bitcoin and gold prices Basher et al. [8]. Illa et al. [9] proposed a methodology for estimating pattern-matching expectations by employing artificial intelligence techniques such as RF and support vector machines.

Artificial intelligence techniques are being more widely used in the stock market due to their capacity to process large amounts of data and identify intricate patterns that humans sometimes find difficult. It is important to keep in mind that these strategies' initial parameter configuration has a significant impact on how effective they are. Inaccurate estimates and outcomes might arise from improperly configured beginning settings. As such, it's important to pay close attention to the parameter settings while using artificial intelligence in stock trading. Consequently, there are numerous optimization algorithms that can be used to overcome these restrictions, such as the whale optimization algorithm (WOA) [10], particle swarm optimization (PSO) [11], Aquila optimizer (AO) [12], battel royal optimization (BRO) [13], biogeography-based optimization (BBO) [14], genetic algorithm (GA) [15], grey wolf optimization (GWO) [16], moth–flame optimization (MFO) [17], and others, can be used to get around these limitations. In 2015, S. Mirjalili proposed the MFO algorithm [17]. The MFO algorithm is a stochastic optimization technique inspired by the natural navigational mechanism of moths. Moths may navigate by keeping their angle concerning a far-off light source, like the moon or a flame, constant. The MFO method leverages this idea to optimize complex problems by varying the position and brightness of synthetic moths, which act as potential solutions. The algorithm explores the issue space to find the optimal answer as fast as possible. Justifications for Selecting the Proposed Model:

Addressing Complex Interactions: The model under consideration adeptly manages the complex interplay between investor psychology, market forces, economic variables, and market dynamics, which poses a significant challenge in the realm of financial forecasting. The comprehensive methodology, which merges the Moth Flame optimization and random forest processes, has been purposefully developed to address the intricacies that are intrinsic to the stock market.

Capability to Adapt to Non-Linear and Non-Stationary Conditions: Predicting stock price time series data presents a significant challenge due to their non-stationary, non-linear, and hypervolatile characteristics. The model being proposed is customized to effectively navigate these intricacies, rendering it well-suited for depicting the ever-changing dynamics of the stock market.

Addressing Computational Complexities: While conventional approaches may improve the accuracy of forecasts, they impose significant computational burdens. The complexities are effectively addressed by the proposed hybrid model, which guarantees precise predictions while maintaining computational efficiency.

Outstanding Efficacy and Performance: In comparison to the alternative models examined in the study, the hybrid model that integrated Moth Flame optimization and the random forest method consistently exhibited superior efficacy and performance. This demonstrates its resilience in addressing the unique difficulties associated with forecasting stock prices.

Optimal Performance and Minimal Error: The proposed model demonstrated an exceptional degree of effectiveness, characterized by minimal error. Preciseness is of the utmost importance when it comes to financial forecasting, as it enables one to make well-informed decisions.

The method proposed is thoroughly elucidated, effectively tackling the complex obstacles inherent in financial forecasting. The complex interaction among economic variables, market forces, and investor psychology poses a substantial obstacle to the ability to make precise forecasts in funding. The intricacy of this matter is compounded by the pronounced volatility, nonstationarity, and non-linearity of stock price time series data within the securities market. Acknowledging the constraints of conventional approaches, the article presents an innovative framework that adeptly surmounts these obstacles through the integration of the Moth Flame optimization methodology and the random forest method. Not only does this hybrid model demonstrate exceptional effectiveness, but it also surpasses other modern models examined in the study. The proposed model exhibits exceptional performance, minimal error, and high efficacy, providing a potentially viable resolution to the computational

intricacies that are intrinsic to conventional forecasting approaches. The research assesses the predictive model that has been proposed by employing Nasdaq index data that covers the period from January 1, 2015, to June 29, 2023. The conclusive findings validate the efficacy and dependability of the suggested model in the domains of stock price analysis and prediction. The experimental results demonstrate that this methodology exhibits a higher level of predictive accuracy in comparison to other approaches. In brief, the methodology that has been proposed effectively tackles the complex issues associated with financial forecasting. It presents an innovative and successful approach that outperforms current models in terms of effectiveness and performance. The reliability of the model is further reinforced through its exhaustive evaluation and validation using real-world data, thereby establishing it as a significant contribution to the domain of stock price prediction. The research investigated many models, including RF, GA-RF, and PSO-RF, to assess their respective levels of reliability. The inquiry encompassed a comprehensive analysis of the data source and all relevant components in the subsequent section. A variety of analytical tools, including optimizer methods, evaluation metrics, and the RF model, were employed to examine the data. The findings of the study are given and compared with those obtained by alternative methodologies in the third part. The fourth part gives information about discussion of the results. The findings of the investigation are succinctly examined in the concluding section.

## II. METHODS AND MATERIALS

### A. Random Forest

A well-liked machine learning technique for situations involving both regression and classification is the Random Forest algorithm, as seen in Fig. 1. It is a subset of the supervised learning algorithms of the support vector machines family. Two other popular tree-based techniques are naive

Bayes and Adaboost. Breiman et al. [6] developed and presented the method, which is well known for being simple and efficient. The RF algorithm creates a variety of intricate decision trees, which improves forecast accuracy. The model is constructed decision trees by selecting the optimal feature from a given collection of characteristics in a non-deterministic manner, resulting in a lower level of predictability compared to alternative tree-based techniques. The methodology operates by iteratively training several decision trees through the utilization of bootstrapping, normalization, and bagging techniques. The grouping strategy involves the simultaneous construction of several decision trees, each using different subsets of characteristics and training data chunks. By employing bootstrapping to ensure the distinctiveness of each decision tree, the variance of the RF is reduced. The RF technique exhibits a high level of promise for generalizability due to the use of many tree-based models for evaluation. By employing this approach, the RF classifier can successfully mitigate challenges associated with unbalanced datasets and overfitting, hence surpassing the performance of current methodologies in accurately recognizing data.

Moreover, the methodology was specifically developed to address the analysis of datasets characterized by a large number of dimensions and strong interdependencies among variables. The reliability of the results increases proportionally with the number of trees included in the ensemble. The high level of precision exhibited by the RF approach can be attributed to the amalgamation of outputs derived from several decision trees. The utilization of ensemble methods mitigates the issue of overfitting and improves the predictive capabilities of the algorithm. Machine learning practitioners highly favor the RF method due to its ability to handle missing data and noisy inputs effectively. The above equation may be utilized to calculate the mean square error for an RF:

$$MSE = \frac{1}{N}\sum_{k=0}^{n} \binom{n}{k} (Fi - Yi)b^2 \qquad (1)$$



Fig. 1. The structure of Random forest.

### B. Particle Swarm Optimization

PSO is an approach that, in order to get optimal results, imitates the cooperative habits of a flock of birds or a school of fish. Even if the exact location of the food supply is unclear at the outset, the swarm will nevertheless follow a set of rules to get there. By working together, the swarm is able to locate the source of nourishment. Swarms of fish or birds will eventually arrive at the near-optimal solution at the same moment. By following these three rules—separation, alignment, and cohesiveness—a bird swarm may efficiently navigate the search space and arrive at the correct solution [18]. Particles undergo separation by moving apart from each other to avoid overcrowding. The particles tend to align with their neighboring particles, resulting in a positional update influenced by the cohesion with said neighbors. Kennedy and Eberhart devised the PSO approach as a means to address optimization challenges. This method draws inspiration from the collective behavior observed in a swarm of particles [11]. PSO technique tends to converge rapidly and necessitates a limited number of parameters, hence reducing computational overhead.

Moreover, the likelihood of encountering a suboptimal local solution is reduced because of the extensive exploration conducted by several particles in quest for an optimal solution. In addition, the algorithm possesses an efficient global search mechanism and does not rely on derivatives. In the PSO, each particle searches a large search space to get the best possible answer. The search process begins with the random generation of candidate solutions, also known as particles, in the search space. Particle velocities and fitness scores are typically computed using a weighted mean of classification accuracy and the number of features in the feature subset. This computation aids in updating the velocity and heading of their trajectories after the initial iteration, and the method is continued until the stopping criterion is reached. The PSO algorithm's particles accelerate and decelerate per the following formula:

$$v_{id}^{t+1} = v_{id}^t + C_1 r_1^t (Pbest_{id}^t - x_{id}^t) + C_2 r_2^t (Gbest_{id}^t - x_{id}^t) \quad (2)$$

The velocity of the ith particle at a given time iteration is denoted as $v_{id}^k$ in a search space with d dimensions. The variables $Pbest_{id}^t$ and $Gbest_{id}^t$ represent the optimal particle and position for each individual and iteration of the ith function. The parameters $C_1$ and $C_2$ are utilized to modify the velocity of particles, whereas $r_1^t$ and $r_2^t$ represent random values within the range of 0 to 1. Furthermore, the particles in the PSO algorithm have the ability to alter their locations by utilizing the equation shown below:

$$x_{id}^{t+1} = x_{id}^t + v_{id}^{t+1} \quad (3)$$

The variable $x_{id}^t$ represents the spatial coordinates of the ith particle at iteration $t$ inside a search space characterized by $d$ dimensions.

### C. Genetic Algorithm

The Genetic algorithm is a computational approach that emulates the mechanism of natural selection in order to address optimization and search problems [15]. Using this approach, a set of potential solutions referred to as individuals is generated. In order to generate novel individuals, these individuals are subsequently exposed to genetic mechanisms such as mutation, recombination, and selection. The assessment process employed in this study is iterative in nature and is repeated through several generations until a viable solution is identified. Consequently, the utilization of GA is prevalent throughout several areas, including but not limited to engineering, finance, and science [19]. GA is comprised of three fundamental components [20]. A chromosome refers to a sequence of numerical or textual symbols that are assigned to each individual by the encoding entity. The selection of an appropriate encoding technique is contingent upon the specific issue that has to be addressed.

Furthermore, the fitness function is utilized to evaluate the degree to which each individual's representation of the answer is accurate. The fitness function has been particularly tailored to address the current problem. To generate novel individuals from existing ones, the evolutionary operators employ the mechanisms of selection, crossover, and mutation. A crossover is a genetic process that combines the chromosomes of two individuals to generate a novel offspring. Mutation, on the other hand, introduces random alterations to an individual's chromosomes. Selection is employed to identify the most reproductively successful individuals.

### D. Mouth Flame Optimization

The Moth Flame Optimizer is a computational model that draws inspiration from natural phenomena and is specifically influenced by the nighttime behavior of butterflies, which is displayed in Fig. 3. [17]. Butterflies have a consistent behavior of fluttering towards the moon when they are attracted by a light source. The Moth Flame Optimizer utilizes and formalizes this approach into an optimization algorithm, which is illustrated in Fig. 2. The optimizer exhibits versatility in its applicability to many optimization issues across several domains, including power and energy systems, economic dispatch, engineering design, image processing, and medical applications. Additionally, the optimizer derives inspiration from the behavior of butterflies as they navigate light sources. Researchers utilize the transverse orientation as a method to investigate the phenomenon of moths maintaining a straight flight trajectory toward the moon [21]. This examination explores the possible use of moths as solutions, which possess the ability to navigate in several dimensions, including $1D, 2D, 3D$, and hyperdimensional space, by manipulating their position vectors. The focus of this study is on examining the spatial distributions of these moths, as they represent the aspects under consideration. The provided methodology guarantees convergence, and the Multi-Objective Firefly Algorithm (MFO) is both reliable and computationally efficient. MFO is commonly utilized as:

$$M = \begin{bmatrix} CO_{1,1} & CO_{1,2} & \cdots & \cdots & CO_{1,h} \\ CO_{2,1} & CO_{2,2} & \cdots & \cdots & CO_{2,h} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ CO_{a,1} & CO_{a,2} & \cdots & \cdots & CO_{n,h} \end{bmatrix} \quad (4)$$

In this context, h represents the number of dimensions, whereas a represents the number of moths.

$$S = \begin{bmatrix} S_{1,1} & S_{1,2} & \cdots & \cdots & S_{1,h} \\ S_{2,1} & S_{2,2} & \cdots & \cdots & S_{2,h} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ S_{a,1} & S_{a,2} & \cdots & \cdots & S_{2,h} \end{bmatrix} \qquad (5)$$

The process of global optimization is conducted by the implementation of the three-step MFO approach.

$$MFO = (I, F, T) \qquad (6)$$

Function $I$ denotes a specific mathematical function, while $F$ symbolizes the flight pattern of a moth as it navigates its environment in search of suitable space. Additionally, the symbol $T$ is used to indicate the criteria that determine when the moth's flight comes to a halt.

$$X_i = t(C_i, S_j) \qquad (7)$$

The formula employed in this context involves the twisting function denoted as $t$, the number of the i-th moths represented by $C_i$, and the number of the $j$-th flames denoted as $S_j$:

$$S(C_i, S_j) = Z_i \cdot e^{bt} \cdot cos(2\pi t) + S_j \qquad (8)$$

The variable $Z_i$ represents the spatial separation between the moth and the flame. The constant $b$ is a parameter in the context of this study. Additionally, the variable t is a random number selected from the interval [-1, 1].

$$Zi = |S_j - X_i| \qquad (9)$$

### E. Data Collection and Preparing

To conduct a comprehensive analysis, it is important to include the trade volume as well as the open, high, low, and closing (OHLC) prices within a certain temporal interval. The data collection period was from January 2, 2015, to June 29, 2023, during which data was obtained from the Nasdaq index on the Yahoo Finance website. The precise details are encompassed inside the dataset that was employed for the investigation. A thorough data-cleaning procedure was conducted to ensure the accuracy and consistency of the forecasting models following the collection of the dataset. The implementation of a multi-step method was devised to safeguard the integrity of the dataset and prevent the inclusion of erroneous or incomplete data that might potentially lead to complications. The data were subjected to a thorough analysis in order to discover any anomalies, extreme numbers, or contradictions that could potentially compromise the validity of the results. This was one of the key aims. Several processes were utilized in order to clean and prepare the data in order to guarantee that it could be utilized. Several procedures, including scaling and normalization, were performed on the data in order to reduce the likelihood of gradient mistakes and unpredictable training results. Before beginning training, the data were normalized by employing the MinMaxScaler method. This was done in order to construct a stable model and reduce the likelihood of extremely high weight values occurring. This normalization process was achieved by employing the equation [22].

$$Xscaled = \frac{(X - Xmin)}{(Xmax - Xmin)} \qquad (10)$$



Fig. 2.   The framework of MFO.

Fig. 3. Movement of the propeller towards the light source.

The training data provided to the model consisted of prices and volume for OHLC. The model was evaluated by including all features except for the close price data. The data were divided into two sets: 80% for training and 20% for testing, as seen in Fig. 4.



Fig. 4. Dividing the data into training and test.

### F. Evaluation Metrics

The accuracy of the projected future was evaluated using a range of performance measures. Carefully selected, these measures provide a comprehensive assessment of the forecasts' validity and accuracy. The assessment method took into account a number of criteria. The Root Mean Square Error (RMSE) determines the root mean square of the errors between the predicted and actual values, the Mean Absolute Percentage Error (MAPE) computes the average absolute difference between the predicted and actual values, and the Coefficient of Determination ($R^2$) quantifies the percentage of the dependent variable's variance that can be predicted based on the independent variable. These techniques help with and are very helpful for evaluating the forecasting models' accuracy.

$$MAPE = \left(\frac{1}{n}\sum_{i=1}^{n}\left|\frac{y_i - \hat{y}_i}{y_i}\right|\right) \times 100 \qquad (11)$$

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \qquad (12)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}} \qquad (13)$$

## III. RESULTS AND DISCUSSION

### A. Hyperparameters Setting

Performance is significantly impacted by the parameters that machine learning algorithms employ. Under these circumstances, it is imperative to guarantee the precise specification of the model parameters. This particular segment provides a comprehensive explanation of the processes involved in hyperparameter configuration. The three optimizers are utilized in order to optimize the parameters of the RF model. In the context of problem-solving, the RF model demonstrates remarkable comfort in tasks that require classification and regression. In prior discussions, we have examined the upper and lower limits of the parameters employed in the configuration of the RF model. In addition to the optimal values derived by the primary optimizer, a detailed dissection of these limits is provided in Table I. Ultimately, the utilization of this data will aid in the determination of the most

effective parameters for the RF model, thereby augmenting its overall performance.

TABLE I. SETTING AND OBTAINING THE OPTIMAL VALUES OF THE HYPERPARAMETERS

| Random Forest | Upper and lower bounds | Best values |
|---|---|---|
| Maximum depth | [10 -100 and None] | 80 |
| Maximum features | [Auto and squared] | auto |
| Minimum samples Leaf | [1 and 4] | 2 |
| Minimum samples split | [2 and 10] | 2 |
| Number of estimators | [200 and 2000] | 500 |

### B. Statistical Values

This phase of the inquiry encompasses Table II, which presents comprehensive statistical data about the dataset. The inclusion of OHLC price and volume statistics in the table enhances the clarity of the data. To comprehensively and accurately evaluate the data, statistical measures such as mean, count, minimum, maximum, standard deviation (Std.), and variations can be employed.

TABLE II. STATISTICS SUMMARY FOR THE DATA SET

| | Open | High | Low | Volume | Close |
|---|---|---|---|---|---|
| Count | 2137 | 2137 | 2137 | 2137 | 2137 |
| Mean | 8744.356 | 8805.287 | 8677.574 | 3143.8 | 8745.821 |
| Std. | 3332.744 | 3362.163 | 3298.311 | 1551.37 | 3332.058 |
| Min | 4218.81 | 4293.22 | 4209.76 | 706.88 | 4266.84 |
| Max | 16120.92 | 16212.23 | 16017.23 | 11621.19 | 16057.44 |
| Variance | 11107186 | 11304139 | 10878852 | 2406747 | 11102609 |

### C. Outcomes of the Models

The primary objective of this study is to identify and assess the most effective hybrid algorithm for the prediction of stock prices. This research is grounded on the establishment of forecasting models and a comprehensive comprehension of the intricate aspects that impact stock market trends. The primary objective is to provide investors and analysts with valuable

information that enables them to make informed and prudent investment choices. The models were processed by using different optimizers and different hyperparameters which were obtained by using those techniques. Table III and Fig. 5, 6 presents a comprehensive examination of the performance of each model. An exhaustive evaluation of the efficacy of each model is also incorporated.

TABLE III. PREDICTED ASSESSMENT RESULTS FOR BENCHMARKING APPROACHES

| | TRAIN SET | | | TEST SET | | |
|---|---|---|---|---|---|---|
| | $R^2$ | RMSE | MAPE | $R^2$ | RMSE | MAPE |
| RF | 0.979 | 423.25 | 4.96 | 0.974 | 254.85 | 1.61 |
| GA-RF | 0.983 | 382.32 | 4.07 | 0.978 | 234.67 | 1.50 |
| PSO-RF | 0.987 | 333.04 | 3.98 | 0.983 | 206.35 | 1.31 |
| MFO-RF | 0.992 | 260.90 | 1.70 | 0.988 | 173.45 | 1.07 |

## IV. DISCUSSION

This study aims to find the best hybrid stock price prediction algorithm. This study relies on predicting models and a deep understanding of stock market dynamics. The goal is to help investors and analysts make smart investment decisions. Fig. 5, 6 and Table III detail in each model's performance. Also included is a thorough model efficacy

review. RMSE, MAPE, and $R^2$ were the three metrics that were utilized in the evaluation of the data analysis, other measures that were utilized included $R^2$ and RMSE. The ability of the aforementioned metrics to give a comprehensive evaluation of the correctness, dependability, and overall efficacy of the analysis is widely acknowledged and widely accepted. Both with and without the utilization of an optimizer, the performance of the RF model has been evaluated by

utilizing the RMSE, MAPE, and $R^2$ criteria. Through the utilization of this approach, it is possible to acquire a more thorough understanding of the performance of the model and to make educated judgments based on the insights that are obtained. After doing an analysis on both the training and test sets, it was observed that the RF model, while it was not utilizing the optimizer, generated $R^2$ values of 0.979 and 0.974 for the training set and the testing set, respectively. While the MAPE values were 4.96 and 1.61, the RMSE values for the training and testing sets were 423.25 and 254.85, respectively. This is in contrast to the MAPE values, which were 4.96 and 1.61. There was a significant improvement in the effectiveness of the RF model as a result of the incorporation of optimizers. An increase in the $R^2$ value to 0.983 for the training dataset and 0.978 for the testing dataset is evidence that the utilization of the GA optimizer has led to considerable improvements. This can be observed by comparing the numbers. The RMSE and the MAPE have both shown a decline in both the training data set and the testing dataset. To be more specific, the RMSE values for the training set were 382.32 and the testing set 234.67, respectively. As an additional point of interest, the MAPE values for the training set are 4.07, whereas the values for the testing set are 1.50. As a result of doing a comparative

analysis between the GA-RF model and the PSO-RF model, it has been concluded that the latter model shows superior performance. For the training phase, the $R^2$ values for the PSO-RF model were found to be 0.987, and for the testing phase, they were found to be 0.983. Despite the fact that the RMSE and MAPE values for training and testing both declined to 333.04 and 3.98, respectively, the values for testing decreased to 206.35 and 1.31, respectively. This is an important observation to make. According to these findings, the PSO-RF model is superior to the GA-RF model in terms of both its degree of effectiveness and its level of efficiency. The remarkable $R^2$ values of 0.992 and 0.988 for training and testing, respectively, demonstrate the efficacy of the MFO-RF model through its remarkable performance. Despite having the lowest possible testing value of 1.07, it performed exceptionally well. The MFO-RF model performed the best when compared to the other models, with the lowest RMSE values of 260.90 for training and 173.45 for testing. For comparison, the other models performed the worst. Considering that these results demonstrate how highly precise and reliable the MFO-RF model is, it is clear that it is an effective instrument for this particular application.

## TRAIN



Fig. 5. The results obtained for $R^2$, RMSE, and MAPE by the proposed model during training.

## TEST



Fig. 6. The results obtained for $R^2$, RMSE, and MAPE by the proposed model during testing.

After thorough research, the MFO-RF model is a reliable instrument for accurately forecasting stock values, thereby establishing its credibility. The efficacy of the model may be assessed by examining the Nasdaq index curves and comparing them to the corresponding curves depicted in Fig. 7 and Fig. 8. The MFO-RF model has superior performance in forecasting

stock prices compared to other models, such as RF, GA-RF, and PSO-RF. A comprehensive analysis of the model's efficacy unveiled that the MFO-RF model predicts stock prices through the integration of the Mouth-Flame optimization technique and the random forest algorithm. The utilization of the RF technique not only mitigates fluctuations in stock prices but also enhances the precision of future trend predictions, hence further augmenting the accuracy of the model. One distinguishing characteristic of the MFO-RF model, in comparison to other models, is its ability to acquire knowledge from previous datasets. To accurately predict stock prices, a model must possess the capability to acquire knowledge from historical datasets and adjust its predictions in response to evolving market patterns. In summary, the MFO-RF model's reliability, accuracy, and ability to derive insights from historical datasets render it a highly valuable tool for predicting stock prices. The utilization of the RF algorithm and MFO optimizer, together with its adaptability in addressing dynamic market trends, positions it as the preferred option for those seeking to achieve profitable stock market transactions. These are the limitations of the research:

Temporal Scope: The study encompasses the period that commences on January 1, 2015, and concludes on June 29, 2023. The temporal scope may fail to encompass specific market conditions or events that transpired beyond the specified time period. Subsequent investigations may delve into the durability of the suggested algorithms across prolonged historical epochs.

Generalization of Algorithms: Although the suggested algorithms demonstrate exceptional performance when applied to Nasdaq data, their ability to be applied to diverse market conditions and financial instrument types is still unknown. It is critical to evaluate the adaptability of the algorithms to a wide range of datasets in order to obtain a thorough comprehension of their practicality.

Insufficient Comparative Research Against Non-Machine Learning Approaches: The study predominantly conducts a comparative analysis of various machine learning models, neglecting to delve deeply into their performance in relation to conventional forecasting methods or statistical methodologies. By incorporating these comparisons, a more comprehensive assessment of the proposed algorithms could be achieved.

The research investigates the impact of hyperparameter selection on sensitivity: To optimize hyperparameters, the study employs genetic algorithms, particle swarm optimization, and Moth Flame optimization. Nevertheless, the explicit consideration of the algorithms' sensitivity to various sets of hyperparameters is absent. A more comprehensive sensitivity analysis may yield valuable insights regarding the models' stability.

Market Volatility Attributable to Inherent Market Volatility: Predicting stock prices necessitates addressing vagaries. The accuracy of the proposed algorithms is commendable; however, the unpredictability inherent in financial markets may introduce unforeseen fluctuations that have an adverse effect on the precision of forecasts.

The potential compromise of model interpretability may arise from the complexity of the proposed algorithms, particularly when they are integrated with optimization techniques. Comprehending the fundamental mechanisms that propel predictions may present a formidable task, thereby constraining the model's applicability in specific decision-making contexts.

Alterations in Market Dynamics: Throughout the period under analysis, the study presupposes a stable set of market dynamics. Variations in economic policies, geopolitical occurrences, or worldwide economic transformations might give rise to modifications in market conduct that the suggested algorithms might not sufficiently account for.

Overfitting Risk: It is crucial to recognize the potential for overfitting, particularly when algorithms are being optimized for particular datasets. While the models might exhibit outstanding performance on the training data, they might encounter difficulties when implemented on unseen data, which raises doubts about their efficacy in real-world scenarios.



Fig. 7. The prediction curve generated via MFO-RF during training.

Fig. 8. The prediction curve generated via MFO-RF during testing.

## V. CONCLUSION

The process of projecting stock prices is a complex and involved task that involves several interconnected components. The stock market is subject to several influences, including but not limited to politics, society, and the economy. It is a complex and always-changing system. To make accurate predictions on future stock values, it is important to consider a range of financial statements, earnings reports, market trends, and other relevant elements. Moreover, it is important to note that macroeconomic factors, such as interest rates, inflation, and worldwide market conditions, wield significant influence over the behavior of the stock market. Developing accurate and dependable prediction models can pose significant challenges owing to the intricate nature and multitude of factors inherent in forecasting stock prices. Understanding the unpredictable and non-linear nature of the market is essential to making accurate forecasts. Fortunately, the MFO-RF model offers a practical answer to these problems and has shown to be accurate and trustworthy. The effectiveness of many stock price prediction models, such as RF, GA-RF, and PSO-RF, was assessed in the current study. By employing the GA, PSO, and MFO hyperparameter optimization techniques, the RF's parameters were enhanced.

Nevertheless, when combined with RF, the MFO optimizer method produced the best results. The OHLC prices and volume for the Nasdaq index from January 2, 2015, to June 29, 2023, made up the dataset utilized in the study. The results of the investigation demonstrate how reliable and accurate the MFO-RF model is at forecasting stock prices.

Throughout the study, the accuracy and predictive capability of the MFO-RF model were evaluated by comparing it to many other models. Based on the obtained data, it can be concluded that the MFO-RF model consistently exhibited superior performance compared to the other models. The testing $R^2$ score of 0.988 indicates a good level of accuracy in the predictions made. The RMSE value of the model, which was observed to be 173.45, indicates that the model's predictions exhibited a satisfactory level of accuracy. The model had a low MAPE score of 1.07, suggesting a consistent ability to provide reliable predictions. In terms of accuracy and efficacy, the MFO-RF model demonstrated superior performance compared to the other models that were examined.

The MFO-RF model is a helpful resource for stock price prediction in general and offers insightful data to investors who are attempting to make well-informed investment decisions.

### REFERENCES

[1] C. Zhang, J. Ding, J. Zhan, and D. Li, "Incomplete three-way multi-attribute group decision making based on adjustable multigranulation Pythagorean fuzzy probabilistic rough sets," International Journal of Approximate Reasoning, vol. 147, pp. 40–59, 2022.

[2] Y.-H. Wang, C.-H. Yeh, H.-W. V. Young, K. Hu, and M.-T. Lo, "On the computational complexity of the empirical mode decomposition algorithm," Physica A: Statistical Mechanics and its Applications, vol. 400, pp. 159–167, 2014, doi: https://doi.org/10.1016/j.physa.2014.01.020.

[3] M. M. Kumbure, C. Lohrmann, P. Luukka, and J. Porras, "Machine learning techniques and data for stock market forecasting: A literature review," Expert Syst Appl, vol. 197, no. December 2021, p. 116659, 2022, doi: 10.1016/j.eswa.2022.116659.

[4] Y. Chen and Y. Hao, "A feature weighted support vector machine and K-nearest neighbor algorithm for stock market indices prediction," Expert Syst Appl, vol. 80, pp. 340–355, 2017.

[5]  A. J. Myles, R. N. Feudale, Y. Liu, N. A. Woody, and S. D. Brown, "An introduction to decision tree modeling," Journal of Chemometrics: A Journal of the Chemometrics Society, vol. 18, no. 6, pp. 275–285, 2004.

[6]  L. Breiman, "Random forests," Mach Learn, vol. 45, pp. 5–32, 2001.

[7]  H. J. Park, Y. Kim, and H. Y. Kim, "Stock market forecasting using a multi-task approach integrating long short-term memory and the random forest framework," Appl Soft Comput, vol. 114, p. 108106, 2022.

[8]  S. A. Basher and P. Sadorsky, "Forecasting Bitcoin price direction with random forests: How important are interest rates, inflation, and market volatility?," Machine Learning with Applications, vol. 9, p. 100355, 2022.

[9]  P. K. Illa, B. Parvathala, and A. K. Sharma, "Stock price prediction methodology using random forest algorithm and support vector machine," Mater Today Proc, vol. 56, pp. 1776–1782, 2022.

[10] S. Mirjalili and A. Lewis, "The whale optimization algorithm," Advances in engineering software, vol. 95, pp. 51–67, 2016.

[11] J. Kennedy and R. Eberhart, "Particle swarm optimization," in Proceedings of ICNN'95-international conference on neural networks, IEEE, 1995, pp. 1942–1948.

[12] L. Abualigah, D. Yousri, M. Abd Elaziz, A. A. Ewees, M. A. A. Al-Qaness, and A. H. Gandomi, "Aquila optimizer: a novel meta-heuristic optimization algorithm," Comput Ind Eng, vol. 157, p. 107250, 2021.

[13] T. Rahkar Farshi, "Battle royale optimization algorithm," Neural Comput Appl, vol. 33, no. 4, pp. 1139–1157, 2021.

[14] D. Simon, "Biogeography-based optimization," IEEE transactions on evolutionary computation, vol. 12, no. 6, pp. 702–713, 2008.

[15] S. Mirjalili, "Genetic Algorithm," in Evolutionary Algorithms and Neural Networks: Theory and Applications, Cham: Springer International Publishing, 2019, pp. 43–55. doi: 10.1007/978-3-319-93025-1_4.

[16] S. Mirjalili, S. M. Mirjalili, and A. Lewis, "Grey wolf optimizer," Advances in engineering software, vol. 69, pp. 46–61, 2014.

[17] S. Mirjalili, "Moth-flame optimization algorithm: A novel nature-inspired heuristic paradigm," Knowl Based Syst, vol. 89, pp. 228–249, 2015.

[18] C. W. Reynolds, "Flocks, herds and schools: A distributed behavioral model," in Proceedings of the 14th annual conference on Computer graphics and interactive techniques, 1987, pp. 25–34.

[19] B. Gülmez and E. Korhan, "COVID-19 vaccine distribution time optimization with Genetic Algorithm," 2022.

[20] E. Alkafaween, A. B. A. Hassanat, and S. Tarawneh, "Improving initial population for genetic algorithm using the multi linear regression based technique (MLRBT)," Communications-Scientific letters of the University of Zilina, vol. 23, no. 1, pp. E1–E10, 2021.

[21] S. Mohammad, A. Laith, A. Hamzeh, A. Mohammad, and A. M. Khasawneh, "Moth–flame optimization algorithm: variants and applications," Neural Comput Appl, vol. 32, no. 14, pp. 9859–9884, 2020.

[22] P. J. M. Ali, R. H. Faraj, E. Koya, P. J. M. Ali, and R. H. Faraj, "Data normalization and standardization: a technical report," Mach Learn Tech Rep, vol. 1, no. 1, pp. 1–6, 2014.

# Improving of Smart Health Houses: Identifying Emotion Recognition using Facial Expression Analysis

Yang SHI, Yanbin BU[*]

School of Media Technology, Communication University of China, Nanjing, Nanjing 211172, China

*Abstract*—**Smart health houses have shown great potential for providing advanced healthcare services and support to individuals. Although various computer vision based approaches have been developed, current facial expression analysis methods still have limitations that need to be addressed. This research paper introduces a facial expression analysis technique for emission recognition based on YOLOv4-based algorithm. The proposed method involves the use of a custom dataset for training, validation, and testing of the model. By overcoming the limitations of existing methods, the proposed technique delivers precise and accurate results in detecting subtle changes in facial expressions. Through several experimental and performance evaluation tasks, we have assessed the efficacy of our proposed method and demonstrated its potential to enhance the accuracy of Smart Health Houses. This study emphasizes the importance of addressing emotional well-being in healthcare. As experimental results shown, the prosed method achieved satisfy accuracy rate and the effectiveness of the YOLOv4 model for emotion detection suggests that emotional intelligence training can be a valuable tool in achieving this goal.**

*Keywords—Smart health houses; computer vision; facial expression; emotion recognition; YOLO*

## I. INTRODUCTION

Smart Health Houses (SHH) are the latest trend in healthcare, which aims to provide seamless and efficient medical services to individuals in the comfort of their homes [1], [2]. These houses utilize advanced technologies to monitor the health status of patients and offer personalized treatment plans [3]. In recent years, there has been a significant increase in the development of Smart Health Houses, and researchers have explored various technologies to improve their effectiveness.

The latest advances in the SHH have shown promising results in improving the quality of medical care [4]. These technologies include wearable devices, sensors, and machine learning algorithms that can detect various physiological signals and track the daily activities of patients [5]–[7]. Recently vision-based methods have been studied by many researchers due to extensive applicability [8]–[11]. Specifically, vision-based systems have attracted many researchers due to their non-invasive nature and ability to detect emotional changes in patients [12], [13]. Computer vision-based systems can analyze facial expressions and provide insights into the emotional state of patients, enabling healthcare providers to offer personalized treatment plans [14],

[15]. Therefore, among these technologies, computer vision-based systems have gained significant attention due to their ability to analyze facial expressions and detect emotions accurately.

Facial expression analysis is one of the most widely studied areas in computer vision-based systems for the SHH [13]–[16]. It has been shown that facial expressions are reliable indicators of emotional changes and can provide valuable information about the patient's mental state [17], [18]. Therefore, researchers have focused on developing accurate facial expression analysis methods to improve the effectiveness of the SHH [15].

Existing methods for emotion recognition can be divided into two categories: conventional methods and deep learning-based methods. Conventional methods include feature extraction and classification algorithms, while deep learning-based methods mainly use convolutional neural networks (CNNs) to learn features automatically from raw data. Deep learning-based methods have shown significant improvements in various tasks [19]–[22], including emotion recognition [23]–[25], due to their ability to automatically learn high-level features from raw data. Despite the recent advances in this field, there are still several limitations and research gaps that need to be addressed. By reviewing of previous studies, existing methods for facial expression analysis have limitations in detecting subtle changes in facial expressions, which can lead to inaccurate results. Therefore, it is necessary to develop more advanced and accurate methods for facial expression analysis to improve the effectiveness of Smart Health Houses.

In this paper, we propose a deep learning-based facial expression analysis method using Yolo-based algorithm. We present a custom dataset for facial expression analysis and use it to train, validate and test our model. Our proposed method addresses the limitations of existing methods and provides accurate results in detecting subtle changes in facial expressions. We evaluate the performance of our proposed method through various experimental and performance evaluation tasks and demonstrate its effectiveness in improving the accuracy of Smart Health Houses.

Our research contributions include identifying the research gap in existing facial expression analysis methods, proposing a deep learning-based method using Yolo-based techniques, and providing experimental and performance evaluations of our proposed method.

The rest of this paper structures as follows, Section II present background of study. Section III reviews the related works. Section IV discusses the material and methods. Section V presents experimental results. Finally, the paper concludes in Section VI.

## II. BACKGROUND OF STUDY

Facial expression analysis for smart health houses is the process of using computer vision and machine learning to detect and analyze facial expressions of individuals in a smart health house environment. To conduct a background study on this topic, this study reviews literature on facial expression recognition and analysis methods for applying facial expression analysis to smart health houses. Furthermore, the Yolo algorithm as effective solution has been selected for emotion recognition in this matter and background of this algorithm is discussed as well.

### A. Facial Expression based Emotion Recognition

Emotion recognition using vision-based facial analysis is a field of research that focuses on the development of algorithms and techniques to automatically detect emotions from facial expressions. This approach is based on the idea that facial expressions are reliable indicators of emotions, and that these expressions can be detected and analyzed using computer vision techniques. The study of emotion recognition using facial analysis has gained increasing attention in recent years due to its potential applications in a variety of fields, such as psychology, marketing, and human-computer interaction. There are various approaches to emotion recognition using vision-based facial analysis, including feature-based methods, holistic methods, and deep learning methods. In following section, the details of each approach s are discussed.

*1) Feature-based emotion recognition*: Feature-based methods extract specific facial features, such as the position and shape of the mouth and eyes, to recognize emotions. The Facial Action Coding System (FACS) is a feature-based approach that identifies specific facial movements associated with emotions. Advantages of this method include the ability to detect subtle facial expressions, while disadvantages include the difficulty of defining features and the reliance on manual coding.

*2) Holistic-based emotion recognition*: Holistic methods analyze the entire face as a whole to detect emotions. These methods include the use of geometric and appearance-based features, such as facial landmarks and texture, to classify emotions. Advantages of this method include the ability to capture the dynamic changes of facial expressions and the ease of implementation. However, the main disadvantage is that they may not be able to capture subtle variations in facial expressions.

*3) Deep learning-based emotion recognition*: Deep learning methods, such as convolutional neural networks (CNNs), have gained popularity in recent years due to their ability to learn complex features from data. These methods involve training large neural networks on large datasets of labeled facial expressions to learn to automatically recognize emotions. Advantages of this method include the ability to automatically learn features, which can reduce the need for manual feature selection and labeling, and the ability to handle complex and high-dimensional data. However, the main disadvantage is that they require large amounts of labeled data to train the networks effectively, which can be expensive and time-consuming to collect.

As literature review indicated in emotion recognition methods, each approach has its advantages and disadvantages. Feature-based methods are good at capturing subtle expressions, while holistic methods can capture dynamic changes. Deep learning methods are highly accurate but require a large amount of labeled data. Therefore, deep based can be investigated because of high efficiency and effectiveness.

### B. YOLO Algorithm

YOLO (You Only Look Once) is a popular object detection algorithm that was first introduced by Joseph Redmon et al. in 2016. The conventional Yolo algorithm architecture is shown in Fig. 1. The YOLO algorithm works by dividing the input image into a grid of cells and predicting bounding boxes and class probabilities for each cell. This allows YOLO to detect multiple objects in an image in real-time with high accuracy.

Over the years, several versions of YOLO have been released in the literature, including: YOLOv1: The original version of YOLO, introduced in 2016, which demonstrated real-time object detection with good accuracy. YOLOv2: Introduced in 2017, this version improved the accuracy of the original algorithm by making changes to the network architecture and adding batch normalization. YOLOv3: Released in 2018, YOLOv3 further improved the accuracy of the algorithm by using a feature pyramid network and introducing new techniques like multi-scale predictions and focal loss. YOLOv4: Introduced in 2020, YOLOv4 is currently the latest version of the YOLO algorithm. This version made significant improvements to the network architecture, including the use of spatial pyramid pooling (SPP), cross-stage partial network (CSP), and Mish activation function. Additionally, YOLOv4 introduced new techniques like data augmentation, label smoothing, and object agnostic NMS, which led to significant improvements in accuracy and speed.

Among the existing Yolo version, the YOLOv4 is considered to be better than its predecessors due to its improved accuracy and speed, as well as its ability to detect smaller objects with higher accuracy. It is a highly effective object detection algorithm that has achieved state-of-the-art performance on several benchmark datasets. Its unique combination of features and optimizations has made it one of the most popular object detection algorithms in the computer vision community.

As shown in Fig. 2(A), feature pyramid network, a backbone network, and several detection heads form the foundation of the YOLOv4 framework. The backbone network is in charge of removing features from the input picture, while the feature pyramid network is utilised to create feature maps at various sizes. Bounding boxes and class probabilities for objects at various scales are predicted using the multiple detection heads.

Fig. 1.    The convention YOLO architecture.



Fig. 2.    The YOLOv4 architecture [26].

The YOLOv4 backbone network is built on a modified version of the CSPDarknet network, a deep neural network that combines the benefits of residual and convolutional networks. The neck network, which features a spatial pyramid pooling (SPP) module that enables the network to collect data at various sizes, comes after the backbone network.

The Path Aggregation Network (PAN), a multi-scale feature fusion module that merges feature maps at several sizes to enhance object recognition accuracy, is improved and used in the feature pyramid network in YOLOv4. The Cross Stage Partial Network (CSP), a feature fusion module that aids in decreasing computation while preserving accuracy, is also a component of the feature pyramid network.

In YOLOv4, the several detection heads are intended to forecast item bounding boxes and class probabilities at various sizes. This is accomplished by employing anchor boxes with various aspect ratios and sizes, which are utilised to recognize objects of various sizes and forms. The detecting heads moreover make use of a brand-new activation technique dubbed Mish, which has been demonstrated to raise deep neural network precision.

## III.    RELATED WORKS

The paper in [27] proposed a method for emotion recognition from facial expressions based on Support Vector Machines (SVM) algorithm. The authors used the JAFFE and CK databases to train the model, which achieved an accuracy rate of over 90%. They also tested the model on real-time video data and found it to be effective for emotion recognition in real-time. The authors suggest that the model can be applied in a variety of applications, such as video conferencing and virtual reality. However, limitations of the study include the use of only one ethnicity, which may limit the model's effectiveness in recognizing emotions in people of other ethnicities. The model may also not be effective in recognizing subtle or complex emotions that are difficult to detect from facial expressions alone.

Bisogni et al [13] explored the impact of deep learning approaches on facial expression recognition (FER) in healthcare industries. The authors compared three different models: CNN, RNN, and a hybrid CNN-RNN model. The hybrid model achieved the highest accuracy rate and was more effective in recognizing subtle emotions such as disgust and contempt. The authors suggest that deep learning approaches

can improve FER in healthcare industries, which can lead to better diagnosis and treatment outcomes. However, the study's limitations include the use of only two datasets and the lack of diversity in the datasets. Further research is needed to validate the effectiveness of deep learning models on a larger and more diverse population.

This paper [28] proposed a video analytics-based facial emotion recognition system for smart buildings. The authors used a dataset of facial expressions to train their machine learning model, which used the Haar Cascade Classifier to detect faces and the Local Binary Patterns Histograms (LBPH) algorithm to recognize facial expressions. The system was tested in a real-world smart building environment, and the authors found that it was able to accurately detect facial expressions and classify them into one of seven emotions. The key features of the system include its ability to operate in real-time and to recognize multiple emotions simultaneously. The authors suggest that this system can be used to improve the well-being and safety of occupants in smart buildings, by identifying and responding to emotional cues. However, limitations of the study include the use of a single dataset and a small sample size for testing. Additionally, the system may not be effective in recognizing subtle or complex emotions that are difficult to detect from facial expressions alone.

Rajavel et al [29] presents an IoT-based smart healthcare video surveillance system using edge computing to analyze facial expressions of patients. The system uses a convolutional neural network (CNN) model to classify facial expressions of pain, discomfort, and distress, and sends alerts to healthcare professionals. The system's key features include its ability to operate in real-time, high accuracy rate, and low power consumption. The findings show that the system achieved a high accuracy rate of 94.3% in detecting facial expressions. However, the system has some limitations, including the need for high-quality video input and limited flexibility in detecting other emotions or expressions. Overall, the proposed system has the potential to enhance healthcare monitoring and improve patient outcomes.

The authors in [30] proposes an IoT-based smart health monitoring system for COVID-19 that utilizes facial expression analysis to monitor the health status of individuals. The method involves capturing facial expressions using a camera and analyzing them using a machine learning algorithm to detect COVID-19 symptoms such as coughing, sneezing, and fever. The system also collects other health data such as heart rate and oxygen levels using wearable devices. The key features of the system include real-time monitoring, early detection of symptoms, and remote monitoring capabilities. The findings suggest that the system can accurately detect COVID-19 symptoms with a high level of accuracy. However, the authors acknowledge some limitations such as the need for further validation studies and the potential for privacy concerns with the use of facial recognition technology.

## IV. MATERIAL AND METHODS

This section presents the details of the proposed method in this study. S mentioned earlier, a Yolo base algorithm is used for emotion recognition. Basically, for this detection, a model is required to generate using a dataset. To generate a Yolo model for this purpose, a custom dataset must be prepared first. This dataset needs to consist of images with labeled emotions (e.g. happy, sad, angry, etc.) and bounding boxes around the faces in each image. The bounding boxes help the Yolo model locate and classify emotions in new images.

### A. Yolo-based usage Justifications

This section intends to justify why the Yolov4 is chosen in this study. The usage of YOLOv4 (You Only Look Once version 4) in this study is primarily focused on its exceptional performance in terms of high accuracy rates compared to other object detection methods. YOLOv4 has gained significant attention and popularity within the computer vision community due to its remarkable ability to detect and locate objects in real-time with remarkable precision.

One of the key justifications for choosing YOLOv4 is its state-of-the-art accuracy, which surpasses many existing object detection algorithms. YOLOv4 achieves this by employing a variety of innovative techniques and architectural enhancements. These advancements include the integration of a more powerful backbone network, feature pyramid network (FPN), spatial pyramid pooling (SPP), and PANet (Path Aggregation Network). These components work collaboratively to extract multi-scale features from the input image, enabling YOLOv4 to detect objects of varying sizes and appearances accurately. Additionally, YOLOv4 utilizes a highly efficient and optimized detection pipeline, enabling it to process images and videos in real-time. By employing a single-pass approach, YOLOv4 eliminates the need for time-consuming region proposal techniques employed by other methods, such as Faster R-CNN. This efficiency is crucial in scenarios where real-time object detection is required, such as autonomous driving, video surveillance, and robotics applications.

Moreover, YOLOv4 incorporates advanced training strategies, including data augmentation techniques, such as mosaic data augmentation and random shape training, as well as optimization methods like focal loss and learning rate scheduling. These strategies contribute to further improving the accuracy of the model. The comprehensive evaluation of YOLOv4 against other state-of-the-art object detection methods has consistently demonstrated its superior performance. It achieves remarkable accuracy rates while maintaining a high detection speed, making it an ideal choice for various applications that demand both accuracy and efficiency.

Fig. 3. Comparison of Yolov4 and other methods in terms of average precision (*AP*) and FPS [30].

As shown in Fig. 3, the graph illustrates comparison between YOLOv4 and other state-of-the-art object detection algorithms. The X-axis represents the Frames Per Second (FPS), which indicates the detection speed of the algorithms, while the Y-axis represents the average precision rate, which signifies the accuracy of object detection. The graph clearly demonstrates that YOLOv4 outperforms the other object detectors in terms of precision rate.

Fig. 3. Comparison of YOLOv4 and other methods in terms of average precision (AP) and FPS, (a) performance comparion base on AP, (b) performance comparion sof base on AP50.

As the average precision rate increases, YOLOv4 consistently achieves higher values compared to its counterparts. This indicates that YOLOv4 is more effective in accurately detecting objects in various scenarios. Several justifications support the superiority of YOLOv4 in terms of precision rate.

Firstly, YOLOv4 adopts an innovative architecture that incorporates advanced techniques such as feature pyramid network (FPN), spatial pyramid pooling (SPP), and path aggregation network (PANet). These components enable YOLOv4 to extract multi-scale features and capture intricate object details, resulting in improved detection accuracy.

Secondly, YOLOv4 utilizes an efficient single-pass detection pipeline, which eliminates the need for time-consuming region proposal techniques. This allows YOLOv4 to process images and videos in real-time without compromising accuracy. Other detectors, such as Faster R-CNN, may achieve higher FPS rates but often at the cost of reduced precision.

Therefore, the graph illustrates that YOLOv4 surpasses other object detectors in terms of precision rate. The innovative architecture, efficient detection pipeline, and advanced training

strategies employed by YOLOv4 contribute to its superiority in accurately detecting objects.

As shown in Fig. 4, the graph provides a comparison of YOLOv4 with other state-of-the-art object detectors, namely Swin Transformer, EfficientDet, SpineDet, YOLOv3, and PP-YOLO. It plots the average precision rate on the Y-axis, representing the accuracy of object detection, against the latency on the X-axis, indicating the time taken for detection. YOLOv4 consistently outperforms the other detectors in terms of precision rate, showcasing its effectiveness in accurately detecting objects across different scenarios. YOLOv4's superiority in precision rate can be attributed to its advanced architecture, which incorporates techniques such as feature pyramid network (FPN), spatial pyramid pooling (SPP), and path aggregation network (PANet). These components enable YOLOv4 to extract multi-scale features and capture intricate object details, resulting in improved detection accuracy. Additionally, YOLOv4's efficient single-pass detection pipeline eliminates the need for time-consuming region proposal techniques, allowing it to process images and videos in real-time without compromising accuracy.



Fig. 4. Comparison of YOLOv4 and others in terms of AP and latency [32].

The effectiveness of YOLOv4 is further justified by its utilization of advanced training strategies. Data augmentation techniques and optimization methods, such as focal loss and learning rate scheduling, enhance the model's ability to generalize and accurately detect objects in diverse conditions. These strategies contribute to YOLOv4's superior precision rate and its ability to surpass other object detectors in terms of accuracy. As results, the graph clearly illustrates YOLOv4's superiority over other object detectors in terms of precision rate. Its advanced architecture, efficient detection pipeline, and advanced training strategies collectively contribute to its effectiveness in accurately detecting objects. YOLOv4's ability to maintain high precision while achieving real-time detection sets it apart from other state-of-the-art detectors, making it a preferred choice for various applications that demand both accuracy and efficiency.

*B. Dataset Preparation*

In this study, we use a dataset from Robloflow universe resource [31]. The Roboflow is a useful resource for preparing a custom dataset for the Yolo model. It allows users to upload

their images and labels, and then provides tools to clean and augment the data. Pre-processing steps, such as resizing and normalizing the images, can improve the accuracy of the Yolo model. Data augmentation techniques, such as random cropping, flipping, and rotation, can increase the size of the dataset and reduce overfitting.

To further improve the accuracy rate of the Yolo model, more advanced data augmentation techniques can be applied. Mix-up augmentation can generate new training samples by blending pairs of images and their labels. Cutout augmentation can randomly remove parts of an image to force the model to focus on other features and using the data augmentation procedure, the total number of images in the dataset 4540.

In nest step, the prepared and enhanced dataset has to be divided into training, validation, and testing sets. The Yolo model is trained using the training set to identify emotions in fresh photos. The validation set is used to fine-tune the model's hyperparameters, including the learning rate and epoch count. The testing set is employed to assess the model's performance on unobserved data. Table I shows the structure of dataset split for training, validation and testing sets.

TABLE I.     DATASET SPLIT FOR TRAINING, VALIDATION AND TESTING SETS

| Training | Validation | Testing |
|---|---|---|
| 85% | 10% | 5% |
| 3859 | 454 | 227 |

The training module is responsible for optimizing the weights of the Yolo model using back propagation and gradient descent. The validation module measures the accuracy of the model on the validation set and adjusts the hyperparameters accordingly. The testing module evaluates the final accuracy of the model on the testing set. All of these modules require careful tuning and monitoring to ensure that the Yolo model is accurately detecting emotions in new images.

*C. Hyperparameter Tunning*

However, based on the obtained results from our experimentation, we set following hyperparameters to generate the YOLOv4-based model,

*1) Learning rate*: Learning rate is the step size at which the model updates its parameters during training. A starting learning rate for YOLOv4 is 0.001.

*2) Batch size*: Batch size is the number of samples processed in a single forward/backward pass. The batch size for YOLOv4 is 64.

*3) Momentum*: Momentum is the parameter that accelerates the gradient descent algorithm in the relevant direction and dampens oscillations. A good momentum value for YOLOv4 is 0.9.

*4) Number of epochs*: The number of epochs is the number of times the entire training dataset is passed through the model during training. A good number of epochs for YOLOv4 is 100.

*5) Anchor boxes*: Anchor boxes are the predefined boxes used to detect objects in YOLOv4. A good number of anchor boxes for emotion recognition is 3-4.

*6) Input size*: The input size of the image affects the detection accuracy and inference time of YOLOv4. A good input size for YOLOv4 is 416*416.

*7) IOU threshold*: IOU threshold is the minimum intersection over union required to consider a detection as true positive. The IOU threshold for YOLOv4 is 0.5.

*8) Confidence threshold*: Confidence threshold is the minimum confidence score required to consider a detection as valid. The confidence threshold for YOLOv4 is 0.25.

*9) NMS threshold*: NMS threshold is the minimum overlap required between two detections to suppress one of them. The NMS threshold for YOLOv4 is 0.45.

*D. Model Generation*

One the dataset is prepared and hyperparameters are tuned, we can generate a model. To generate this model, we use a pre-trained model, by downloading a pre-trained weight from the COCO dataset. This weight is used to initialize the YOLOv4 model. When a pre-trained weight is available, we train the model using a Yolov4 model, and then validate the model to ensure it's not overfitting, and finally test the model to evaluate its performance.

## V. EXPERIMENTAL RESULTS

This section presents the experimental result and performance evaluation of the proposed method. Experimental results obtained using YOLOv4 model on image samples have shown promising results in terms of object detection accuracy. It has also shown excellent performance in detecting small objects and improving accuracy over the facial emotion recognition dataset. Fig. 5 shows some samples of experimental results.

The created YOLOv4 model has demonstrated outstanding generalization skills, i.e., it can identify emissions in a variety of complicated scenarios, as demonstrated by the experimental findings. Overall, YOLOv4's efficacy and superiority over other object identification models have been shown by experimental findings on picture samples utilizing this model.

For performance evaluation average precision are calculated for by classes. This calculation is performed for validation and testing sets individually. Table II presents the Average precision by class for validation and testing sets. Moreover, precision, recall and Mean Average Precision (*mAP*) metrics are calculated using to measure the performance and corresponding diagram are demonstrated.

Precision is the percentage of accurately identified emotions (true positive predictions) among all projected emotions. True positives divided by the total of true positives and false positives is how it is determined. A high precision means that most of the model's predictions are accurate. Fig. 6 represents the result of precision metric.

Recall quantifies the share of accurate predictions that were positive out of all the actual emotions in the dataset. By dividing the total of true positives and false negatives, it is

calculated. The majority of the emotions in the dataset can be detected by the model, as indicated by a high recall. Fig. 7 represents the result of recall metric.

The mAP scores calculated for each emotion class. AP measures how well the model can distinguish between different emotions. It is calculated as the area under the precision-recall curve for a specific emotion class. A high mAP indicates that the model can accurately detect all the emotion classes. Fig. 8 represents the result of mAP metric.

Therefore, in order to evaluate the efficiency of the model, we use these metrics to evaluate the performance of the models. As experimental results and performance evaluations reported a higher precision and recall indicate better performance, while a higher mAP indicates better differentiation between different emotion classes. By using these metrics, we address the generated model achieve satisfy accuracy rate and the effectiveness of the YOLOv4 model for emotion detection.



Fig. 5. Samples of experimental results.

TABLE II. AVERAGE PRECISION BY CLASS FOR VALIDATION AND TESTING SETS

| Classes | Validation set | Test set |
|---|---|---|
| Angry | 93% | 91% |
| Happy | 87% | 89% |
| Sad | 89% | 85% |
| Surprised | 96% | 94% |
| All | 91.25% | 89.75% |

As shown in Fig. 6, the mAP_0.5 and mAP_0.5:0.95 are metrics used to evaluate the efficiency and effectiveness of object detection models, including the generated YOLOv4 model for facial expression analysis and emotion recognition. mAP_0.5 measures the average precision at an IoU threshold of 0.5, indicating how well the model localizes objects. Higher mAP_0.5 scores indicate better accuracy. mAP_0.5:0.95 considers a range of IoU thresholds and calculates the average precision across this range, providing a comprehensive evaluation of the model's performance. A higher mAP_0.5:0.95 score indicates accurate detection across various IoU thresholds. The mAP curves visualize the precision-recall trade-off, indicating the model's ability to achieve high precision while maintaining a reasonable recall rate. High mAP scores justify accurate results in emotion recognition. The YOLOv4 model demonstrates effectiveness by achieving high mAP scores, indicating accurate detection and recognition of facial expressions for emotion recognition tasks.

As experimental results indicated, this study significantly advances our understanding of emotion recognition through facial expression analysis. Extensive comparisons have been meticulously carried out to demonstrate the superiority of the proposed method over existing approaches, as visually depicted in Fig. 3 and Fig. 4. Moreover, the discussions closely accompany these figures to provide detailed insights and interpretations.

Furthermore, the presentation of the proposed method's performance serves as a pivotal aspect of our study, showcasing the outcomes in a manner that emphasizes their significance. As obtained results and illustrated in Fig. 3 and Fig. 4, the developed method not only exhibits effectiveness but also outperforms other existing methods, underscoring its importance in the field of emotion recognition.

Fig. 6. Performance evaluation based on precision metric.



Fig. 7. Performance evaluation based on recall metric.





Fig. 8. Performance evaluation based on mAP metrics.

## VI. CONCLUSION

Smart Health Houses aim to provide efficient and personalized medical services to individuals in their homes using advanced technologies, such as wearable devices, sensors, and machine learning algorithms. Among these technologies, computer vision-based systems that can analyze facial expressions and detect emotions accurately have gained significant attention due to their non-invasive nature. However, existing methods for facial expression analysis have limitations in detecting subtle changes in facial expressions. To address this issue, this paper proposes a deep learning-based facial expression analysis method using a YOLOv4-based algorithm. The method utilizes a custom dataset for training, validation, and testing and overcomes the limitations of existing methods, providing accurate results. The study concludes that the proposed method has the potential to enhance the accuracy of Smart Health Houses. For directions and future studies, investigating the potential of combining multiple technologies, including computer vision-based systems, wearable devices, and sensors, to enhance the accuracy and effectiveness of Smart Health Houses. This study holds considerable significance as it not only advances our understanding of emotion recognition through facial expression analysis but also offers a superior method for achieving accurate results. The outcomes presented in this research provide a valuable foundation for future studies in the realm of emotion recognition, paving the way for more sophisticated and precise techniques. Researchers can build upon these findings to develop enhanced algorithms and applications, ultimately contributing to a deeper comprehension of human emotions and their applications in various fields, such as psychology, human-computer interaction, and affective computing. Further exploring the limitations of existing facial expression analysis methods and developing more advanced techniques is to improve their accuracy. This could involve investigating the potential of using different deep learning architectures, such as convolutional neural networks, to enhance the performance of facial expression analysis models. Additionally, research could focus on developing methods to address issues such as variations in lighting conditions and occlusions of facial features.

REFERENCES

[1] T. M. Ghazal et al., "IoT for smart cities: Machine learning approaches in smart healthcare—A review," Future Internet, vol. 13, no. 8, p. 218, 2021.

[2] M. M. Islam, A. Rahaman, and M. R. Islam, "Development of smart healthcare monitoring system in IoT environment," SN Comput Sci, vol. 1, pp. 1–11, 2020.

[3] V. Bhardwaj, R. Joshi, and A. M. Gaur, "IoT-based smart health monitoring system for COVID-19," SN Comput Sci, vol. 3, no. 2, p. 137, 2022.

[4] A. Das Gupta, S. M. Rafi, B. R. Rajagopal, T. Milton, and S. G. Hymlin, "Comparative analysis of internet of things (IoT) in supporting the health care professionals towards smart health research using correlation analysis," Bull. Env. Pharmacol. Life Sci., Spl, no. 1, pp. 701–708, 2022.

[5] J. Gao et al., "Ultra‐robust and extensible fibrous mechanical sensors for wearable smart healthcare," Advanced Materials, vol. 34, no. 20, p. 2107511, 2022.

[6] A. Sujith, G. S. Sajja, V. Mahalakshmi, S. Nuhmani, and B. Prasanalakshmi, "Systematic review of smart health monitoring using deep learning and Artificial intelligence," Neuroscience Informatics, vol. 2, no. 3, p. 100028, 2022.

[7] A. A. Nancy, D. Ravindran, P. M. D. Raj Vincent, K. Srinivasan, and D. Gutierrez Reina, "Iot-cloud-based smart healthcare monitoring system for heart disease prediction via deep learning," Electronics (Basel), vol. 11, no. 15, p. 2292, 2022.

[8] A. Aghamohammadi, M. C. Ang, E. A. Sundararajan, K. W. Ng, M. Mogharrebi, and S. Y. Banihashem, "Correction: A parallel spatiotemporal saliency and discriminative online learning method for visual target tracking in aerial videos," PLoS One, vol. 13, no. 3, p. e0195418, 2018.

[9] V. Nandini, R. D. Vishal, C. A. Prakash, and S. Aishwarya, "A review on applications of machine vision systems in industries," Indian J Sci Technol, vol. 9, no. 48, pp. 1–5, 2016.

[10] C.-Z. Dong and F. N. Catbas, "A review of computer vision–based structural health monitoring at local and global levels," Struct Health Monit, vol. 20, no. 2, pp. 692–743, 2021.

[11] Y. Jiang, W. Wang, and C. Zhao, "A machine vision-based realtime anomaly detection method for industrial products using deep learning," in 2019 Chinese Automation Congress (CAC), IEEE, 2019, pp. 4842–4847.

[12] S. P. Yadav, S. Zaidi, A. Mishra, and V. Yadav, "Survey on machine learning in speech emotion recognition and vision systems using a recurrent neural network (RNN)," Archives of Computational Methods in Engineering, vol. 29, no. 3, pp. 1753–1770, 2022.

[13] C. Bisogni, A. Castiglione, S. Hossain, F. Narducci, and S. Umer, "Impact of deep learning approaches on facial expression recognition in healthcare industries," IEEE Trans Industr Inform, vol. 18, no. 8, pp. 5619–5627, 2022.

[14] P. Silapasuphakornwong and K. Uehira, "Smart mirror for elderly emotion monitoring," in 2021 IEEE 3rd Global Conference on Life Sciences and Technologies (LifeTech), IEEE, 2021, pp. 356–359.

[15] F. A. Pujol, H. Mora, and A. Martínez, "Emotion recognition to improve e-healthcare systems in smart cities," in Research & Innovation Forum 2019: Technology, Innovation, Education, and their Social Impact 1, Springer, 2019, pp. 245–254.

[16] J. K. Pandey, V. Veeraiah, S. B. Talukdar, V. B. Talukdar, V. M. Rathod, and D. Dhabliya, "Smart City Approaches Using Machine Learning and the IoT," in Handbook of Research on Data-Driven Mathematical Modeling in Smart Cities, IGI Global, 2023, pp. 345–362.

[17] A. P. Plageras and K. E. Psannis, "IOT-based health and emotion care system," ICT Express, vol. 9, no. 1, pp. 112–115, 2023.

[18] R. Jahangir, Y. W. Teh, F. Hanif, and G. Mujtaba, "Deep learning approaches for speech emotion recognition: State of the art and research challenges," Multimed Tools Appl, pp. 1–68, 2021.

[19] A. Aghamohammadi, R. Ranjbarzadeh, F. Naiemi, M. Mogharrebi, S. Dorosti, and M. Bendechache, "TPCNN: two-path convolutional neural network for tumor and liver segmentation in CT images using a novel encoding approach," Expert Syst Appl, vol. 183, p. 115406, 2021.

[20] W. Mellouk and W. Handouzi, "Facial emotion recognition using deep learning: review and insights," Procedia Comput Sci, vol. 175, pp. 689–694, 2020.

[21] A. M. Abu Nada, E. Alajrami, A. A. Al-Saqqa, and S. S. Abu-Naser, "Age and gender prediction and validation through single user images using cnn," 2020.

[22] R. Ranjbarzadeh et al., "Lung infection segmentation for COVID-19 pneumonia based on a cascade convolutional network from CT images," Biomed Res Int, vol. 2021, pp. 1–16, 2021.

[23] N. Mehendale, "Facial emotion recognition using convolutional neural networks (FERC)," SN Appl Sci, vol. 2, no. 3, p. 446, 2020.

[24] D. K. Jain, P. Shamsolmoali, and P. Sehdev, "Extended deep neural network for facial emotion recognition," Pattern Recognit Lett, vol. 120, pp. 69–74, 2019.

[25] H. Ge, Z. Zhu, Y. Dai, B. Wang, and X. Wu, "Facial expression recognition based on deep learning," Comput Methods Programs Biomed, vol. 215, p. 106621, 2022.

[26] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," arXiv preprint arXiv:2004.10934, 2020.

[27] P. Tarnowski, M. Kołodziej, A. Majkowski, and R. J. Rak, "Emotion recognition using facial expressions," Procedia Comput Sci, vol. 108, pp. 1175–1184, 2017.

[28] K. S. Gautam and S. K. Thangavel, "Video analytics-based facial emotion recognition system for smart buildings," International Journal of Computers and Applications, vol. 43, no. 9, pp. 858–867, 2021.

[29] R. Rajavel, S. K. Ravichandran, K. Harimoorthy, P. Nagappan, and K. R. Gobichettipalayam, "IoT-based smart healthcare video surveillance system using edge computing," J Ambient Intell Humaniz Comput, pp. 1–13, 2022.

[30] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Scaled-yolov4: Scaling cross stage partial network," in Proceedings of the IEEE/cvf conference on computer vision and pattern recognition, 2021, pp. 13029–13038.

[31] "Facial Emotion Rrecognition Dataset," Roboflow Universe, 2023.

[32] M. Murugavel, "YOLO V4.," https://manivannan-ai.medium.com/yolo-v4-750cd627064f.

# Perceived Benefits and Challenges of Implementing CMMI on Agile Project Management: A Systematic Literature Review

Anggia Astridita, Teguh Raharjo, Anita Nur Fitriani

Faculty of Computer Science, University of Indonesia, Jakarta, Indonesia

*Abstract*—**In an era where the agility and responsiveness of Agile project management are paramount, the integration of structured models like the Capability Maturity Model Integration (CMMI) presents a blend of unique opportunities and challenges. This study conducts a comprehensive systematic literature review of 23 scientific articles, chosen through the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) methodology, to explore the benefits and challenges of CMMI and software development integration within the context of Agile project management. Emphasizing the enhancement of Agile project management maturity, the research delves into the role of CMMI, particularly CMMI-DEV, as a pivotal element in Software Process Improvement (SPI) models tailored to Agile environments. The study's novelty lies in its systematic and in-depth investigation of CMMI's integration with Agile project management methodologies, a critical yet underexplored area in the existing literature. Addressing the urgency highlighted by global trends of resource inefficiencies and project management challenges, this research offers timely insights for both academia and industry. This study also categorizes key benefits while identifying prevalent challenges, such as resource constraints and organizational resistance. Additionally, this research also suggests solutions and improvements to these challenges. By offering a comprehensive evaluation, the research significantly advances the understanding of the complexities and potential of CMMI and Agile project management integration. It provides valuable insights for practical applications in organizational settings, emphasizing the potential of integrating structured models like CMMI-DEV with Agile project management methodologies. This integration is essential for enhancing project management maturity, marking a significant step forward in academic research and practical applications in this vital domain.**

*Keywords*—*CMMI; SPI; Agile project management; systematic literature review; PRISMA*

## I. INTRODUCTION

In today's rapidly evolving business landscape, the agility and efficiency of project management have become pivotal for organizational success [1]. Agile project management, in particular, has emerged as a critical determinant in enabling companies across diverse sectors to navigate the complexities of fluctuating market demands [2]–[4]. However, the transition to agile methodologies is not without its challenges. Research and industry observations underscore significant inefficiencies in current project management practices, negatively impacting organizational performance [4].

The 2023 Pulse of the Profession survey by the Project Management Institute reveals a striking indicator of these challenges. This survey reports a global average loss of 5.2% in investment due to subpar project performance, marking a sustained trend of resource wastage over recent years [5]. Such findings signal a more profound, systemic issue in project management across various industries, calling for a strategic approach to enhance project management maturity [4].

In response to this need, frameworks like the Project Management Maturity Model (PMMM) have been developed, offering a structured way to evaluate and uplift an organization's project management practices [4]. Yet, in the face of the intricate challenges posed by the modern business environment, achieving maturity in project management alone is insufficient. The significance of Software Process Improvement (SPI) and frameworks like the Capability Maturity Model Integration (CMMI) becomes more evident in this context. They provide comprehensive methodologies to bolster an organization's software development and management prowess, crucial in the digital era [6], [7].

The CMMI model has evolved beyond its initial focus on software engineering to support organizations across various industries in enhancing their capabilities, measuring performance, and addressing common business challenges. The CMMI Model represents a set of established best practices that can be applied globally to help organizations develop key capabilities [15]. It is designed to be user-friendly, adaptable, and compatible with other methodologies, such as Agile, SAFe, and DevSecOps, among others [8], [14], [15]. However, integrating CMMI with other methodology such as Agile project management is not straightforward, often requiring significant investment and grappling with the stringent requirements of these frameworks [8], [9]. This indicates that the integration process is complex and resource intensive.

Despite these hurdles, many organizations, especially those prioritizing high-quality outputs, are increasingly exploring the synergies between agile methods and CMMI [8], [10], [11]. The latest CMMI V2.0 has been recognized for its improvements in project management performance [4], [12], [13], as acknowledged by professionals worldwide [6]. Therefore, it is essential to examine both the benefits and challenges of integrating CMMI with Agile project management to fully understand the value of such a combination.

Previous research has examined the relationship between the challenges of combining CMMI and Agile. Henriquez et al. [8] conducted a study to determine how much CMMI addresses these challenges. They focused on two significant CMMI artifacts for integration and emphasized vital issues that organizations must address [8]. Additionally, Ferdinansyah et al. [10] compiled experiences from combining software and explored the challenges involved in the collaborative implementation process. Their research also delved into the compatibility between CMMI and Agile Development [10]. Henriquez et al. [14] also conducted another research that focuses on analyzing and identifying agile artifacts that align with CMMI-DEV V2.0 practices, aiding Agile organizations in adopting or transitioning to this latest model from CMMI-DEV V1.3, with a particular emphasis on Planning and Managing Work Practice Areas and Practices. However, it is worth noting that previous studies have not explicitly examined the Perceived Benefits and Challenges of Implementing CMMI in Agile project management.

This research analyzes the potential benefits and challenges associated, that comes from the integration of CMMI-DEV and Agile project management. The aim is to leverage the advantages of implementing CMMI and Agile project management. Moreover, the study will also identify the challenges and recommend solutions accordingly. Both academic research and organizational practice can benefit from this research. The proposed solutions for these challenges can serve as valuable guidance for senior managers considering the joint implementation of CMMI and Agile project management. Additionally, it provides the most recent literature review, which can be utilized to enhance research on CMMI and agile project management in academic research. This study aims to address the subsequent research questions:

RQ1: What are the benefits of integrating CMMI and Agile Project Management?

RQ2: What are the challenges and the corresponding solution of integrating CMMI and Agile Project Management?

## II. LITERATURE REVIEW

### A. Agile Project Management

Agile project management (APM) methodology, a methodology that was developed approximately two decades ago predominantly for the software sector, draws its principles from the Agile Manifesto. These principles prioritize individual interactions, the functionality of software, cooperation with customers, and flexibility in adapting to change [3]. Originally for software development, APM is now utilized in various fields, characterized by team autonomy, iterative development, and equality within teams [1], [2]. It aims to deliver high-value products within time and budget constraints by integrating planning with execution and fostering teamwork and customer collaboration. APM manages paradoxical dynamics, balancing team flexibility with procedural rigor, and is recognized for effectively responding to evolving project requirements and customer needs [16].

Research and practitioner experiences indicate that APM positively impacts behavioral, affective, and cognitive outcomes, with a more pronounced effect on behavioral aspects like performance and innovation. This effectiveness is not limited to software development; APM shows slightly greater effectiveness in non-software domains [3]. The adaptability and broader application of APM highlight its role in enhancing project management practices and fostering positive changes in workplace behavior and performance. These insights reflect APM's comprehensive impact across different industries. The pervasiveness of APM's influence underscores its potential to revolutionize project management practices and drive organizational success across diverse domains [1], [3].

### B. Software Process Improvement (SPI)

It is vital to enhance both the efficiency and effectiveness of software development and management processes. In this regard, the role of Software Process Improvement (SPI) is pivotal for accomplishing such improvements [7], [17]. SPI involves developing and honing a set of collective knowledge and practices specific to software development, with a continuous commitment to improving these processes [6]. This ongoing improvement helps organizations increase their development performance and adapt efficiently to evolving business environments, thereby gaining a competitive edge. SPI is a critical enabler for ensuring that software development and management processes are aligned with business goals and objectives [17], [18].

A range of SPI (Software Process Improvement) models are employed within the software industry. These include the Capability Maturity Model (CMM), Capability Maturity Model Integration (CMMI), People Software Process (PSP), SPI, Capability Determination (SPICE), and BOOTSTRAP. [6]. These models address various challenges, including projects exceeding budgets and timelines, subpar software quality, and unmet requirements. Different SPI models, including CMMI, PSP, SPICE, MSF, RUP, ISO, IDEAL, and Six Sigma, are employed to tackle these issues [6], [19]. The CMMI model is particularly notable for its wide-ranging application across different sectors, extending beyond software. It helps organizations improve their processes, leading to better quality and efficiency in software development projects [6], [7], [20].

### C. CMMI-DEV

CMMI was developed as a process improvement model to help many organizations improve performance, achieve process maturity, and achieve organizational goals [14], [21]. The CMMI model encompasses three different categories. CMMI-DEV focuses on product and service development; CMMI-SVC is dedicated to establishing and managing services; and CMMI-ACQ pertains to acquiring products and services. This paper explicitly centers on CMMI-DEV, the model used in computer programming, underscoring its two essential cycle regions for necessities [22], [23].

Additionally, the CMMI Institute recently introduced version 2.0 of the CMMI. With its emphasis on continuous improvement and process optimization, CMMI-DEV V2.0 empowers organizations to enhance their software development capabilities and deliver exceptional products that meet stakeholder needs [14], [24]. This updated version integrates practices from the three version 1.3 constellations (DEV, ACQ, and SVC) and the People Capability Maturity Model (PCMM). In CMMI V2.0, CMMI for Supplier

Management (CMMI-SPM) will replace CMMI-ACQ from V1.3 [25]. The mapping between CMMI 1.3 and CMMI 2.0 models is shown in Fig. 1.

Moreover, the latest version of CMMI introduces several new practice areas, with project management being a key area predominantly addressed in CMMI-DEV [24]. The practice areas consist of Estimating (EST), Planning (PLAN), Risk and Opportunity Management (RSK), Monitor and Control (MC), Implementation Infrastructure (II), Requirements Development and Management (RDM), Supplier Agreement Management (SAM) [26]. Therefore, this paper focuses on the new version of CMMI, particularly CMMI-DEV.



Fig. 1. CMMI 1.3 and CMMI 2.0 models mapping.

### D. Project Management Category of CMMI

In CMMI, several process areas target the project management domain. The project management encompasses several areas, such as project planning, monitoring and control, and many others [9], [27]. This paper primarily centered on project management, specifically the management activities associated with project planning, monitoring, and control. Table I maps process areas in CMMI 1.3 to practice areas in CMMI 2.0 relevant to project management based on its definition for each process areas accordingly.

TABLE I. CMMI PROCESS AREAS AND PRACTICE AREA MAPPING RELATING TO PROJECT MANAGEMENT

| Process Area CMMI 1.3 [20] | Process Area CMMI 2.0 [26] |
|---|---|
| Project Planning (PP) | Estimating (EST), Planning (PLAN), Risk & Opportunity Management (RSK) |
| Project Monitoring and Control (PMC) | Monitor & Control (MC), RSK |
| Integrated Project Management (IPM) | PLAN, Implementation Infrastructure (II), MC |
| Risk Management (RSKM) | RSK |
| Requirement Management (REQM) | Requirements Development & Management (RDM) |
| Supplier Agreement Management (SAM) | Supplier Agreement Management (SAM), PLAN, MC |

## III. RESEARCH METHOD

Systematic Literature Reviews (SLRs) are commonly employed to conclude, gather evidence, and produce concise summaries. In this research, the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) methodology is employed [28]. PRISMA is known for its transparent and concise approach, making it preferable over other methods. While initially developed for healthcare

research, PRISMA has also proven effective in Information System studies [28]. The PRISMA Workflow for this research is depicted in Fig. 1, providing a visual representation of the process.

### A. Planning the SLR

In this stage, keyword generation is conducted based on the main keywords obtained at the beginning of the research, namely CMMI and Agile Project Management. The keywords ("CMMI" AND "Agile" AND "Project Management") were identified. These keywords were used to construct queries on each database. The query formats were subsequently modified to align with the specific query requirements on the advanced search function of each database.

The online database with extensive collections of SPI papers relevant to the problem domain and criteria of the studies was selected based on specific considerations. The reason is due to the relevant search results available, the appropriate field usually provided by the online database, and the reputation of the online database itself. For this research, the selected databases include the ACM Digital Library, IEEE Xplore, Scopus, and Science Direct. The criteria applied in the study selection process are detailed in Table II.

TABLE II. CRITERIA FOR THE STUDY SELECTION

| | | |
|---|---|---|
| **Inclusion Criteria (IC)** | IC1 | The research was written in English. |
| | IC2 | The research was published between 2018 and 2023. |
| | IC3 | The research in the computer science problem domain. |
| | IC4 | The research is academic research, such as a journal or conference paper. |
| **Exclusion Criteria (EC)** | EC1 | The research was not relevant to CMMI implementation in Agile settings. |
| | EC2 | The research does not mention the challenges or benefits of CMMI implementation in Agile settings. |

### B. Implementation of SLR

In this stage of the SLR, a meticulous and expansive search was conducted to identify pertinent literature. Sections relevant to the formulated research questions underwent thorough analysis. The analysis entailed an in-depth examination of the entire text, specifically emphasizing investigations related to combining CMMI and Agile Project Management methodologies.

Fig. 2 delineates the SLR workflow using PRISMA, illustrating the systematic selection process. This process commenced with the initial data acquisition from various databases, adhering to the inclusion and exclusion criteria outlined in Table I, and was completed with the identification of results that corresponded with the predefined research criteria. The identified studies were then subjected to rigorous quality assessment to ensure the validity and reliability of the findings.

A total of 23 scientific articles, sourced from a range of highly ranked journals and conference proceedings, as per Scimago, have been selected for in-depth analysis regarding the challenges and effects on project management.

Fig. 2.   Workflow of literature review (PRISMA).

## C. Reporting the SLR

The final stage of the SLR involves the presentation of findings. First, Table III is a comprehensive reference guide, summarizing all relevant articles aligned with the research goals. This table includes article titles, publication years, journal index information, and reference citation coding.



Fig. 3.   Distribution of scientific articles per publication year.

Secondly, Fig. 3 visually represents scientific papers categorized according to their publication years. Table III and Fig. 3 enhance the presentation and comprehension of the SLR's outcomes, ensuring a robust and comprehensive exploration of the research area.

TABLE III.      SUMMARY OF ALL RELEVANT ARTICLES

| No | Title | Year | Index | Code |
|----|-------|------|-------|------|
| 1 | Agile-CMMI Alignment: Contributions and To-Dos for Organizations | 2021 | Q2 | [8] |
| 2 | Bringing to Light the Agile Artifacts Pointed Out by CMMI | 2022 | Q1 | [14] |
| 3 | A model for defining project lifecycle phases: Implementation of CMMI level 2 specific practice | 2022 | Q1 | [29] |
| 4 | Software project management in high maturity: A systematic literature mapping | 2019 | Q1 | [30] |
| 5 | Does agile methodology fit all characteristics of software projects? Review and analysis | 2023 | Q1 | [31] |
| 6 | Software Requirement Analysis: Research Challenges and Technical Approaches | 2018 | Procd | [23] |
| 7 | Abandonment of a Software Process Improvement Program: Insights from Case Studies | 2020 | Procd | [32] |
| 8 | A Conceptual View for an Enhanced Cloud Software Life-Cycle Process (CSLCP) Model | 2020 | Procd | [33] |
| 9 | Reconciliation of scrum and the project management process of the ISO/IEC 29110 standard-Entry profile—an experimental evaluation through usability measures | 2021 | Q2 | [34] |
| 10 | Practical Suggestions to Successfully Adopt the CMMI V2.0 Development for Better Process, Performance, and Products | 2023 | Procd | [12] |
| 11 | Crafting a CMMI V2 Compliant Process for Governance Practice Area: An Experiential Proposal | 2020 | Procd | [24] |
| 12 | Challenges in Combining Agile Development and CMMI: A Systematic Literature Review | 2021 | Procd | [10] |
| 13 | A new scrum and CMMI level 2 compatible model for small software firms in order to enhance their software quality | 2022 | Procd | [11] |
| 14 | The CMMI-Dev Implementation Factors for Software Quality Improvement: A Case of XYZ Corporation | 2020 | Procd | [35] |
| 15 | The Improvement Process for The Software Development and Requirements Management to Achieve Capability Level 3 of CMMI | 2022 | Procd | [36] |
| 16 | Improving the Quality of Requirements Engineering Process in Software Development with Agile Methods: A Case Study Telemedicine Startup XYZ | 2021 | Procd | [37] |
| 17 | Model-driven gap analysis for the fulfillment of quality standards in software development processes | 2023 | Q2 | [38] |
| 18 | Software quality models: Exploratory review | 2023 | Q3 | [39] |
| 19 | Agile Governance Guidelines for Software Development SMEs | 2021 | Procd | [27] |
| 20 | A Novel model to adapt CMMI Level 2 by Assessing the Local SMEs of Bangladesh | 2023 | Procd | [40] |
| 21 | Towards Implementation of Process and Product Quality Assurance Process Area for Saudi Arabian Small and Medium Sized Software Development Organizations | 2018 | Q2 | [41] |
| 22 | Evaluation of Maturity Level of the Electronic based Government System in the Department of Industry and Commerce of Banjar Regency | 2020 | Q3 | [42] |
| 23 | Software Process Improvement During the Last Decade: A Theoretical Mapping and Future Avenues | 2021 | Procd | [43] |

## IV. RESULTS AND DISCUSSION

The eligible studies underwent a process of screening and analysis to extract relevant information. While screening the complete text, the problem domain and research questions were identified concurrently.

### A. Perceived Benefits of CMMI and Agile Project Management Integration

In exploring the symbiosis between CMMI and Agile Project Management, it's essential to delve into the tangible benefits this integration brings to project management. This chapter categorizes these benefits into five core practice areas within CMMI: Estimating, Planning, Risk and Opportunity Management, Monitor and Control, and Requirements Development and Management.

*1) Estimation:* Enhanced Accuracy and Predictability. Studies have shown that integrating CMMI with Agile methodologies can significantly refine budget and schedule predictions. Alqadri et al. [35], Saputra et al. [42], and Galvan-Cruz et al. [34] highlight CMMI's effectiveness in this regard. Degerli, M. [12], [24] citing CMMI Institute [13], reports a remarkable 17% increase in estimation accuracy following the adoption of CMMI V2.0. These enhanced estimation processes reduce project uncertainties and bolster the likelihood of successful project execution. Importantly, Albuquerque [32] and Itzik et al. [31] note the role of CMMI in providing greater predictability in costs and deadlines, contributing to reduced costs and increased productivity.

*2) Planning:* Streamlined Project Management. The incorporation of CMMI into Agile methodologies elevates the process of work planning and management. As outlined by Valeria et al. [8], [35], CMMI's structure enables the creation of comprehensive forecasts concerning workload, costs, and schedules. This foresight is crucial in preventing budget or timeline overruns [8], [14]. The synergy of Agile and CMMI also improves goal attainment. It minimizes rework, as evidenced by the significant reduction in rework (70%) and the increase in on-time delivery rates (97%) reported by Degerli, M. [12], [24] referencing the CMMI Institute [13].

*3) Risk:* Mitigating Risks with Informed Strategies. In Agile environments, effectively handling requirement changes is critical. CMMI's quality models play a significant role here, as they aid in risk reduction and quality enhancement [23], [37]. The model promotes a proactive approach to identifying and evaluating risks and opportunities, as noted by Degerli, M. [12], [24]. This approach encompasses establishing performance benchmarks derived from historical data, facilitating early identification of variances, and supporting informed choices in project management.

*4) Monitor and control:* Ensuring Quality and Compliance. CMMI's role in monitoring and controlling project and organizational processes is pivotal [10]. It provides a framework for analyzing and managing critical subprocesses, particularly in high-maturity project management scenarios, as discussed by Cerdeiral, C. T., & Santos, G. [30] Keshta et al. [41]. This supervisory function ensures that software quality is upheld during its development and that the end products fulfill users' expectations.

*5) Requirements development and management:* Optimizing Product Quality and Customer Satisfaction. A vital aspect of CMMI is its guidance on requirement development and management [23]. This element holds particular significance in Agile Project Management, given its focus on flexibility and adaptability. According to several studies, CMMI enhances customer satisfaction by improving product quality and aligning closely with customer needs [29], [31], [35], [36]. Furthermore, CMMI's principles can be seamlessly amalgamated with Agile's emphasis on customer collaboration and adaptability to change [36], [39].

### B. Classification of Challenges of CMMI and Agile Project Management Integration

The eligible studies highlighted several challenges faced by organizations. A thematic categorization focusing on tactical and organizational challenges was employed, as delineated by Valeria et al. [8].

*1) Resource and Time Constraints:* Integrating CMMI into Agile projects presents significant resource and time challenges. Ferdinansyah et al. [10] highlight that this integration demands additional resources, effort, and time beyond the scope of standard project activities. Adopting new concepts and practices within an Agile framework requires careful planning and considerable investment. These constraints are particularly impacting Small and Medium-sized Enterprises (SMEs), as observed by Henríquez et al. [27] and Saheel et al. [11]. SMEs often operate with limited budgets and resources, making it challenging to sustain the additional demands of integrating CMMI [33], [38], [40], [41], [43]. This challenge can lead to difficulties in maintaining the balance between the pursuit of quality improvement and the practical realities of project management within these organizations.

*2) Organizational Resistance and Change Management:* Resistance to adopting new methodologies is a common challenge within organizations. This resistance often stems from a lack of knowledge or experience with the new systems, as noted by Valeria et al. [8], Ferdinansyah et al. [10], Albuquerque et al. [32], and Demirel & Das [23]. Frequent leadership changes or past experiences with unsuccessful management initiatives can further exacerbate such resistance. This skepticism and disinterest, especially among long-standing employees, can pose significant hurdles to successfully integrating CMMI and Agile methodologies. Moreover, the lack of support from top management in enforcing new processes can lead to demotivation among practitioners and quality assurance teams, hampering the overall adoption process [8], [10], [23], [32].

*3) Balancing Agility with Control:* A critical challenge in integrating CMMI with Agile methodologies is balancing the structured approach of CMMI with the flexibility of Agile [8].

As organizations strive to achieve higher maturity levels within the CMMI framework, they may find that agility is compromised, as pointed out by Valeria et al. [8]. Ferdinansyah et al. [10] further elaborate that the control and accountability emphasized by CMMI can clash with the core principles of Agile, which values adaptability and minimal bureaucratic overhead. This challenge concerns managing project processes and aligning organizational culture and values with these contrasting methodologies.

*4) Knowledge, Training, and Expertise Gaps:* The successful adoption of CMMI in an Agile setting relies heavily on sufficient knowledge, training, and expertise. Valeria et al. [8] and Albuquerque et al. [32] highlight the challenges organizations face due to a lack of in-depth understanding of CMMI and Agile methodologies. The deficiency in specialized training covering the full spectrum of development activities can hinder the effective implementation of these methodologies. Furthermore, the absence of support from specialists in statistical and process management knowledge can leave project managers and process groups struggling to align organizational and project goals with critical subprocesses [30], [38].

*5) Scaling and Knowledge Dissemination:* Scaling and disseminating practices aligning with CMMI and Agile across an organization is a significant challenge, especially at higher maturity levels. Valeria et al. [8] emphasized that without the active support of upper management, experiences, and practices beneficial to integrating CMMI and Agile often remain confined to specific teams or projects. This limitation prevents these practices from being adopted more widely throughout the organization. Scaling experiences effectively requires documentation of successful practices and a concerted effort to share and institutionalize these practices across various teams and departments [8], [30].

### C. Challenges and Solution Mapping

This paper aims to align the solutions by drawing on the identified challenges. Table IV presents the mapping of the challenges and their corresponding potential solutions. This structured approach facilitates a clearer understanding of how each solution directly addresses the specific challenges, thereby enhancing the overall effectiveness of the proposed solutions. Additionally, this alignment ensures that the proposed solutions effectively address the identified challenges, paving the way for a more robust and impactful implementation strategy.

TABLE IV.     SOLUTION MAPPING

| No | Challenges | Solution |
|---|---|---|
| 1 | Resource and Time Constraints | Senior managers must decisively commit to process improvement initiatives, defining the scope and allocating the right resources. The success of these efforts often hinges on continuous investment in process improvements and addressing barriers such as resource limitations, inexperienced staff, organizational politics, and time pressure. The CMMI model, increasing popularity, can be a crucial tool for such improvements. However, understanding how to sustain these improvements post-appraisal, especially under time and budget constraints, remains vital for long-term success. Both higher and lower-level management's support and commitment play a critical role in this regard [8], [12], [27]. |
| 2 | Organizational Resistance and Change Management | Addressing human issues like resistance and acceptance is essential for senior managers. Strong management support is vital in implementing changes, especially adopting an Agile philosophy. This support helps facilitate cultural change and overcome initial resistance from factors like lack of human resources and work overload. Agile teams may show resistance even with management support, indicating the need for a more comprehensive approach to change management [8], [27], [32]. |
| 3 | Balancing control and agility | Effective integration of CMMI with Agile processes requires tools that automate control and accountability. This approach helps align traditional governance processes with Agile teams and addresses the additional work CMMI might introduce. Constant monitoring of agility is crucial during this alignment. Developing prescriptive guidelines for Process Areas aligned with business goals can help manage the impact on skill [8], [10], [27], [36]. |
| 4 | Knowledge, Training, and Expertise Gaps | Senior management must play an active role in contextualizing governance and providing organizational training for Agile practices. Addressing knowledge and experience gaps in CMMI and Agile Development is critical, especially for higher CMMI Maturity Levels. The lack of training for new employees in process improvement is a significant issue that needs to be addressed to bridge knowledge and expertise gaps. Practices from CMMI V2.0 and SAFe 5.0 can support this process [8], [27], [36]. |
| 5 | Scaling and Knowledge Dissemination | Senior managers should adopt Agile strategies that decentralize decision-making and organize work based on business value. Implementing quality standards and tools for verifying development processes, alongside a model-driven approach, can enable a comprehensive assessment of software development. Addressing Agile's scaling, training, and organizational policy limitations is critical. Knowledge dissemination poses a challenge in large organizations, necessitating the development of non-bureaucratic processes that meet organizational needs. Software Process Improvement (SPI) methodology can extend Agile to enhance product innovation, quality, and efficiency [8], [27], [36], [38]. |

## V. CONCLUSIONS

This paper has explored the perceived benefits and challenges of this integration. Integrating CMMI with Agile Project Management presents a promising avenue for enhancing project management practices. Based on further analysis, this integration can lead to more accurate estimations, streamlined project planning, effective risk mitigation, improved quality control, and optimized requirement development and management. However, this integration comes with challenges such as resource constraints, organizational resistance to change, and the need to balance control and agility. Bridging knowledge gaps and scaling integrated practices across the organization are notable hurdles.

Senior management must play a pivotal role by committing to process improvement initiatives, providing necessary resources, and supporting cultural change to overcome challenges. Automation tools for control and accountability, along with continuous monitoring of agility, are essential for

maintaining the balance between CMMI and Agile. Additionally, addressing knowledge gaps through training and promoting non-bureaucratic knowledge dissemination processes are vital steps towards successful integration. Overall, a strategic and committed approach from senior management is crucial to realizing the full potential of the CMMI and Agile integration and improving project management practices for better project outcomes.

### A. Limitations of Study

The study's primary limitation lies in integrating CMMI-DEV and software development within Agile project management. The literature review's scope, constrained by the time frame (2018-2023) and language of the publications (English), may also limit the comprehensiveness of the findings.

### B. Future Works

Future research should focus on empirical studies to validate the findings of this systematic literature review. Investigations involving case studies or surveys in various organizational settings would offer deeper insights into the practical implementation challenges and benefits of combining CMMI-DEV and software development within Agile project management. Additionally, exploring the evolution of these frameworks in rapidly changing technological landscapes will provide more dynamic and current insights into their integration and application.

REFERENCES

[1] P. Putri, T. Raharjo, B. Hardian, and T. Simanungkalit, "Challenges and Best Practices Solution of Agile Project Management in Public Sector: A Systematic Literature Review," JOIV : International Journal on Informatics Visualization, vol. 7, Nov. 2023, doi: 10.30630/joiv.7.2.1098.

[2] R. M. Haj Hamad and M. Al Fayoumi, "Scalable Agile Transformation Process (SATP) to Convert Waterfall Project Management Office into Agile Project Management Office," in 2018 International Arab Conference on Information Technology (ACIT), 2018, pp. 1–8. doi: 10.1109/ACIT.2018.8672701.

[3] J. Koch, I. Drazic, and C. C. Schermuly, "The affective, behavioural and cognitive outcomes of agile project management: A preliminary meta-analysis," J Occup Organ Psychol, vol. 96, no. 3, pp. 678–706, Sep. 2023, doi: 10.1111/joop.12429.

[4] E. Fabbro and S. Tonchia, "Project Management Maturity Models: Literature Review and New Developments," The Journal of Modern Project Management, vol. 8, no. 3, Apr. 2022, doi: 10.19255/JMPM02503.

[5] Project Management Institute, "Pulse of the Profession® 2023, 14th Edition," 2023. Accessed: Nov. 10, 2023. [Online]. Available: https://www.pmi.org/learning/thought-leadership/pulse/power-skills-redefining-project-success

[6] A. Singh and S. S. Gill, "Measuring the maturity of Indian small and medium enterprises for unofficial readiness for capability maturity model integration-based software process improvement," Journal of Software: Evolution and Process, vol. 32, no. 9, Sep. 2020, doi: 10.1002/smr.2261.

[7] C. Y. Chen and J. C. Lee, "Comparative effects of knowledge-based antecedents in different realms of CMMI-based software process improvement success," Comput Stand Interfaces, vol. 81, Apr. 2022, doi: 10.1016/j.csi.2021.103599.

[8] V. Henriquez, A. M. Moreno, J. A. Calvo-Manzano, and T. S. Feliu, "Agile-CMMI Alignment: Contributions and To-Dos for Organizations," Computer (Long Beach Calif), vol. 54, no. 12, pp. 38–49, Dec. 2021, doi: 10.1109/MC.2020.3039105.

[9] A. B. Farid and Y. M. Helmy, "Implementing Project Management Category Process Areas of CMMI Version 1.3 Using Scrum Practices, and Assets," IJACSA) International Journal of Advanced Computer Science and Applications, vol. 7, no. 2, 2016, Accessed: Aug. 11, 2023. [Online]. Available: https://thesai.org/Publications/ViewPaper?Volume=7&Issue=2&Code=IJACSA&SerialNo=34

[10] A. Ferdinansyah and B. Purwandari, "Challenges in Combining Agile Development and CMMI: A Systematic Literature Review," in ACM International Conference Proceeding Series, Association for Computing Machinery, Feb. 2021, pp. 63–69. doi: 10.1145/3457784.3457803.

[11] S. Saheel, L. Bin Rahman, F. Humaira, F. Sadia, and M. Hasan, "A new scrum and CMMI level 2 compatible model for small software firms in order to enhance their software quality," in Proceedings of the ACM Symposium on Applied Computing, Association for Computing Machinery, Apr. 2022, pp. 1607–1610. doi: 10.1145/3477314.3507173.

[12] M. Degerli, "Practical Suggestions to Successfully Adopt the CMMI V2.0 Development for Better Process, Performance, and Products," in 2020 5th International Conference on Computer Science and Engineering (UBMK), 2020, pp. 126–129. doi: 10.1109/UBMK50275.2020.9219438.

[13] CMMI Institute, "Take Organizational Performance to the Next Level with CMMI," 2021. Accessed: Nov. 10, 2023. [Online]. Available: https://www.isaca.org/resources/infographics/take-organizational-performance-to-the-next-level-with-cmmi

[14] V. Henriquez, J. A. Calvo-Manzano, A. M. Moreno, and T. San Feliu, "Agile-CMMI V2.0 alignment: Bringing to light the agile artifacts pointed out by CMMI," Comput Stand Interfaces, vol. 82, p. 103610, Aug. 2022, doi: 10.1016/J.CSI.2021.103610.

[15] ISACA, "CMMI by the Numbers: How Companies Worldwide Are Improving Performance," Sep. 2023.

[16] D. Binci, C. Cerruti, G. Masili, and C. Paternoster, "Ambidexterity and Agile project management: an empirical framework," TQM Journal, vol. 35, no. 5, pp. 1275–1309, Jun. 2023, doi: 10.1108/TQM-01-2022-0011.

[17] S. Badshah, A. A. Khan, and B. Khan, "Towards Process Improvement in DevOps: A Systematic Literature Review," in ACM International Conference Proceeding Series, Association for Computing Machinery, Apr. 2020, pp. 427–433. doi: 10.1145/3383219.3383280.

[18] A. Al-Ashmori, P. D. D. Dominic, S. Basri, A. Muneer, and G. Naji, "Literature Review: Blockchain-Oriented Software Characteristics and New Stream for Software Process Improvement," in 2022 International Conference on Decision Aid Sciences and Applications, DASA 2022, Institute of Electrical and Electronics Engineers Inc., 2022, pp. 905–910. doi: 10.1109/DASA54658.2022.9765124.

[19] M. Ahmad and Z. A. Rana, "Comparative Analysis of light weight practices for SPI in small and medium software organizations," in 2021 15th International Conference on Open Source Systems and Technologies (ICOSST), 2021, pp. 1–6. doi: 10.1109/ICOSST53930.2021.9683961.

[20] A. Hidayati, B. Purwandari, E. K. Budiardjo, and I. Solichah, "Global Software Development and Capability Maturity Model Integration: A Systematic Literature Review," in 2018 Third International Conference on Informatics and Computing (ICIC), 2018, pp. 1–6. doi: 10.1109/IAC.2018.8780489.

[21] CMMI Institute, A Guide to Scrum and CMMI ®: Improving Agile Performance with CMMI 2. 2016.

[22] D. Proença and J. Borbinha, "Formalizing ISO/IEC 15504-5 and SEI CMMI v1.3 – Enabling automatic inference of maturity and capability levels," Comput Stand Interfaces, vol. 60, pp. 13–25, Nov. 2018, doi: 10.1016/j.csi.2018.04.007.

[23] S. Demirel and R. Das, "Software Requirement Analysis: Research Challenges and Technical Approaches," 2018, Accessed: Aug. 11, 2023. [Online]. Available: https://ieeexplore.ieee.org/document/8355322

[24] M. Degerli, "Crafting a CMMI V2 Compliant Process for Governance Practice Area: An Experiential Proposal," in 2020 Turkish National Software Engineering Symposium (UYMS), 2020, pp. 1–3. doi: 10.1109/UYMS50627.2020.9247068.

[25] ISACA, "What is the difference between CMMI, CMMI-DEV, CMMI-SVC, CMMI-SPM (formerly CMMI-ACQ), People CMM, and DMM?"

Accessed: Nov. 10, 2023. [Online]. Available: https://support.isaca.org/s/article/What-is-the-difference-between-CMMI-CMMI-DEV-CMMI-SVC-CMMI-SPM-formerly-CMMI-ACQ-People-CMM-and-DMM-1598331745508

[26] CMMI Institute, "CMMI V2.0 to V1.3 Practice Mapping," 2018. Accessed: Nov. 09, 2023. [Online]. Available: https://cmmiinstitute.com/resource-files/public/v2-0-materials/cmmi-v2-0-to-v1-3-practice-mapping-(%E6%B1%89%E8%AF%AD)

[27] V. Henriquez and A. M. Moreno, "Agile Governance Guidelines for Software Development SMEs," in Iberian Conference on Information Systems and Technologies, CISTI, IEEE Computer Society, Jun. 2021. doi: 10.23919/CISTI52073.2021.9476224.

[28] M. J. Page et al., "The PRISMA 2020 statement: an updated guideline for reporting systematic reviews," Syst Rev, vol. 10, no. 1, p. 89, 2021, doi: 10.1186/s13643-021-01626-4.

[29] I. Keshta, "A model for defining project lifecycle phases: Implementation of CMMI level 2 specific practice," Journal of King Saud University - Computer and Information Sciences, vol. 34, no. 2, pp. 398–407, Feb. 2022, doi: 10.1016/j.jksuci.2019.10.013.

[30] C. T. Cerdeiral and G. Santos, "Software project management in high maturity: A systematic literature mapping," Journal of Systems and Software, vol. 148, pp. 56–87, Feb. 2019, doi: 10.1016/j.jss.2018.10.002.

[31] D. Itzik and G. Roy, "Does agile methodology fit all characteristics of software projects? Review and analysis," Empir Softw Eng, vol. 28, no. 4, Jul. 2023, doi: 10.1007/s10664-023-10334-7.

[32] R. Albuquerque, G. Santos, A. Malucelli, and S. Reinehr, "Abandonment of a Software Process Improvement Program: Insights from Case Studies," in ACM International Conference Proceeding Series, Association for Computing Machinery, Dec. 2020. doi: 10.1145/3439961.3439966.

[33] A. A. Alshazly, M. Y. Elnainay, A. A. El-Zoghabi, and M. S. Abougabal, "A Conceptual View for an Enhanced Cloud Software Life-Cycle Process (CSLCP) Model," in ACM International Conference Proceeding Series, Association for Computing Machinery, Nov. 2020, pp. 1–5. doi: 10.1145/3436829.3436830.

[34] S. Galvan-Cruz, M. Mora, C. Y. Laporte, and H. Duran-Limon, "Reconciliation of scrum and the project management process of the ISO/IEC 29110 standard-Entry profile—an experimental evaluation through usability measures," Software Quality Journal, vol. 29, no. 2, pp. 239–273, Jun. 2021, doi: 10.1007/s11219-021-09552-3.

[35] Y. Alqadri, E. K. Budiardjo, A. Ferdinansyah, and M. F. Rokhman, "The CMMI-Dev Implementation Factors for Software Quality Improvement: A Case of XYZ Corporation," in ACM International Conference Proceeding Series, Association for Computing Machinery, Jan. 2020, pp. 34–40. doi: 10.1145/3379310.3379327.

[36] P. Joembunthanaphong and G. Sriharee, "The Improvement Process for The Software Development and Requirements Management to Achieve Capability Level 3 of CMMI," in ICSEC 2022 - International Computer Science and Engineering Conference 2022, Institute of Electrical and Electronics Engineers Inc., 2022, pp. 296–301. doi: 10.1109/ICSEC56337.2022.10049356.

[37] A. S. Wibawa, E. K. Budiardjo, and K. Mahatma, "Improving the Quality of Requirements Engineering Process in Software Development with Agile Methods: A Case Study Telemedicine Startup XYZ," in 2021 International Conference Advancement in Data Science, E-Learning and Information Systems, ICADEIS 2021, Institute of Electrical and Electronics Engineers Inc., 2021. doi: 10.1109/ICADEIS52521.2021.9701962.

[38] G. Giachetti, J. L. de la Vara, and B. Marín, "Model-driven gap analysis for the fulfillment of quality standards in software development processes," Software Quality Journal, 2023, doi: 10.1007/s11219-023-09649-x.

[39] L. Pinedo et al., "Software quality models: Exploratory review," EAI Endorsed Transactions on Scalable Information Systems, vol. 10, no. 6, 2023, doi: 10.4108/eetsis.3982.

[40] F. Oyshi et al., "A Novel model to adapt CMMI Level 2 by Assessing the Local SMEs of Bangladesh," in Procedia Computer Science, Elsevier B.V., 2023, pp. 2043–2050. doi: 10.1016/j.procs.2023.01.506.

[41] I. Keshta, M. Niazi, and M. Alshayeb, "Towards Implementation of Process and Product Quality Assurance Process Area for Saudi Arabian Small and Medium Sized Software Development Organizations," IEEE Access, vol. 6, pp. 41643–41675, Jul. 2018, doi: 10.1109/ACCESS.2018.2859249.

[42] M. R. Y. Saputra, W. W. Winarno, Henderi, and S. Shaddiq, "Evaluation of maturity level of the electronic based government system in the department of industry and commerce of banjar regency," Journal of Robotics and Control (JRC), vol. 1, no. 5, pp. 156–161, Sep. 2020, doi: 10.18196/jrc.1532.

[43] A. Al-Ashmori, P. D. D. Dominic, S. Basri, Q. Al-Tashi, A. Muneer, and E. A. A. Ghaleb, "Software Process Improvement during the Last Decade: A Theoretical Mapping and Future Avenues," in 2021 International Congress of Advanced Technology and Engineering, ICOTEN 2021, Institute of Electrical and Electronics Engineers Inc., Jul. 2021. doi: 10.1109/ICOTEN52080.2021.9493426.

# Crime Prediction Model using Three Classification Techniques: Random Forest, Logistic Regression, and LightGBM

Abdulrahman Alsubayhin, Muhammad Sher Ramzan, Bander Alzahrani

King Abdulaziz University, Jeddah, Kingdom of Saudi Arabia

*Abstract*—**Predicting the likelihood of a crime occurring is difficult, but machine learning can be used to develop models that can do so. Random forest, logistic regression, and LightGBM are three well-known classification methods that can be applied to crime prediction. Random forest is an ensemble learning algorithm that predicts by combining multiple decision trees. It is an effective method for classification tasks, and it is frequently employed for crime prediction because it handles imbalanced datasets well. Logistic regression is a linear model that can be used to predict the probability of a binary outcome, such as the occurrence of a crime. It is a relatively straightforward technique that can be effective for crime prediction if the features are carefully chosen. LightGBM is a gradient-boosting decision tree algorithm with a reputation for speed and precision. It is a relatively new algorithm, but because it can achieve high accuracy on even small datasets, it has rapidly gained popularity for crime prediction. The experimental results show that the LightGBM performs best for binary classification, followed by Random Forest and Logistic Regression.**

*Keywords—Crime prediction; random forest; logistic regression; LightGBM*

## I. INTRODUCTION

In a culture where crime is low, it is always disturbing to see the number of crimes rising [1]. Crime is a social issue that hinders the economic growth of a nation. Crime has always existed, and violent crime is the greatest threat to society [2]. Population growth and urbanization have dramatically increased criminal activity [3], particularly in urban areas [4].

In recent years, crime prediction has acquired popularity because it enables investigation authorities to handle crimes computationally [5]. Better predictive algorithms that direct police patrols towards criminals are required [6]. Several research investigations have been conducted to predict crime categories, crime rates, and crime hotspots using crime datasets from various regions, such as South Korea and the United States [7]. Additionally, using the Canada dataset, various prototype projects are expanded to identify crime-related geographic locations, such as residential and commercial areas [8].

Crime threatens us and society, necessitating serious consideration if we expect to reduce its onset or consequences.

Daily, data officers working alongside law enforcement authorities throughout the United States record hundreds of crimes. Numerous cities in the United States have signed the Open Data initiative, making crime data and other categories of data accessible to the public. This initiative aims to increase citizen participation in decision-making by uncovering interesting and valuable facts using this data [9].

San Francisco is one of many cities that have joined this Open Data movement. The data scientists and engineers working with the San Francisco Police Department (SFPD) have documented over one hundred thousand criminal cases based on police complaints [10]. Using these historical data, numerous patterns can be uncovered. This would help us identify crimes that may occur in the future, allowing the municipal police to better protect the city's population [11].

Violent and nonviolent crimes are predicted and classified using random forest, logistic regression, and LightGBM. The primary objective of this paper is to propose a crime prediction model based on past criminal records.

Using three techniques, the proposed model evaluates accuracy, log loss, ROC AUC, precision, and recall evaluation matrices. The data is descriptively analyzed, and crime statistics' spatial and temporal distribution are visualized to identify potential patterns. The original dataset's features are extracted, and classification is carried out using random forest, logistic regression, and LightGBM techniques.

LightGBM has the highest performance for binary classification, followed by random forest and logistic regression, according to the experimental results. LightGBM has the best precision, accuracy, log loss, ROC AUC, and F1 score. It has the lowest recall, but this is not inherently a negative attribute. In this case, the dataset is imbalanced, as there are far more examples of class 0 than class 1. This means that avoiding false positives is essential to avoiding false negatives. LightGBM accomplishes this by emphasizing recall while maintaining high precision. Random forest has a lower accuracy, log loss, ROC AUC, precision, and F1 score than LightGBM but a higher recall. This indicates that random forest is superior at avoiding false negatives but less effective at predicting true positives. Logistic regression has the three models' lowest accuracy, log loss, ROC AUC, precision, and F1 score. It has the lowest recall as well. This indicates that logistic regression is the model with the worst performance for binary classification.

Overall, the model with the greatest performance for binary classification is LightGBM, followed by random forest and logistic regression.

## A. Related Work

Due to the relationship between crime and society, predictions of future crime have been investigated extensively. These studies use machine learning algorithms to address these predictions. Using machine learning algorithms to predict spatial crime data has proven effective [12].

Accurate crime prediction is difficult but essential for preventing criminal behavior. Accurately estimating the crime rate, types, and hot areas based on historical patterns presents numerous computational challenges and opportunities [5].

Prediction analysis is dominated by crime prediction based on machine learning; however, few studies systematically compare machine learning methods. The ability of machine learning algorithms to process non-linear rational data has been validated in numerous disciplines, including crime prediction. It can process high-dimensional data with a faster training pace and extract the characteristics of the data [13].

Despite extensive research efforts, the literature lacks relative accuracy for crime prediction from large datasets for multiple locations, such as Los Angeles and Chicago datasets.

The authors of [14] employ the model to improve the effectiveness of criminal investigation systems. This model identifies crime patterns based on inferences gathered from the crime site and predicts the perpetrator's description of the suspect most likely responsible for the crime. This work has two primary elements: Analyzing and forecasting the perpetrator's identity. The crime analysis phase identifies the number of unsolved crimes and evaluates the impact of variables such as year, month, and weapon on those crimes. The prognosis phase estimates the perpetrators' characteristics, such as age, gender, and relationship to the victim. These hypotheses are based on the evidence gathered at the crime scene. The system predicts the perpetrator's physical characteristics using algorithms such as multilinear regression, K-neighbors classifier, and neural networks. It was trained and evaluated using the San Francisco Homicide dataset (1981-2014) and Python.

Yao et al. used the San Francisco Dataset; this paper is based on the random forest algorithm, which splits the study areas into four groups based on the hot spot distribution based on historical crime data: frequent hot areas, common hot areas, occasional hot areas, and non-hot areas; then, corresponding covariates from non-historical crime data are added to the prediction model to investigate changes in the result accuracy of crime prediction [15]. The data relies on actual data, and the experimental findings reveal that compared to the inference approach based solely on historical crime data, the model with covariates outperforms the model without covariates.

A preliminary analysis of the spatiotemporal crime patterns in San Francisco is attempted in this study [16]. They use spectral analysis to examine the temporal evolution of all crime categories, discovering that many exhibit a weekly or monthly pattern and other components. They demonstrate that spatial distribution has weekly patterns. These findings can be used to develop predictive models for policing and increase knowledge of crime dynamics.

## II. DATA ANALYSIS

The model in the study is built using a Kaggle dataset [17]. The dataset (training set/data) has several properties, each with its own link. The Kaggle incidences of San Francisco crimes are included in the training dataset. The data spans the years January 2003 to May 2015. The collection covers nearly 12 years of San Francisco criminal reports. The collection contains categories of all crimes containing various crime types.

The original training dataset is arbitrarily mixed and divided into training and testing datasets of 80% and 20%, respectively, in the study. Any data imbalances relating to the "Primary Type" feature were corrected using a combination of oversampling (SMOTE) and random sampling; SMOTE stands for Synthetic Minority Over-sampling Technique. It is a data augmentation approach used in machine learning to deal with skewed datasets. SMOTE generates synthetic minority class samples by combining existing minority class samples. This balances class distribution and improves machine learning model performance on minority class predictions.

### A. Features

Every entry in our data set pertains to a specific crime, and each data record includes the following characteristics:

- Dates - The date and time of the crime.

- Category - The type of crime. In the classification stage, we must forecast this target/label.

- Descript - A brief description of any relevant details of the crime.

- DayOfWeek - The weekday on which the offense happened.

- PdDistrict - The Police Department District to which the offense has been assigned.

- Resolution - How the crime was resolved (for example, by arresting or booking the culprit).

- Address - The crime scene's approximate street address.

- X - Longitude of a crime's location.

- Y - The latitude of a crime's site.

### B. Preprocessing

We execute various preprocessing processes on our datasets to achieve better classification results before deploying any algorithms. These are some examples:

- Eliminating features like resolution, description, and address. The resolution and description of a crime are only known after the crime has occurred and have limited relevance in a practical, real-world scenario where one is attempting to predict what type of crime has occurred; hence, these were deleted. We deleted the address since we already knew the latitude and longitude; in that context, the address added little marginal value.

- The weekdays, police, and criminal categories were all indexed and replaced with numbers.

- The timestamp included the year, date, and time of each offense. This was broken down into five components: Year (2003-2015), Month (1-12), Date (1-31), Hour (0-23), and Minute (0-59).

- We used Python's get_dummies() function to turn categorical variables into dummy variables. Dummy variables are binary variables that show whether a given category exists. For example, if a category variable contains three types, the get_dummies() function will generate two dummy variables. One dummy variable indicates the presence of the first category, while the other indicates the presence of the second category. The lack of the first two dummy variables will implicitly reflect the third category.

- In machine learning, the get_dummies() function is a popular technique to deal with categorical variables. Many machine learning algorithms cannot deal with categorical data directly. Machine learning algorithms may train models using data that includes categorical variables by turning categorical variables into dummy variables. We also remove certain unnecessary elements, such as incidentNum and coordinate.

Feature inclusion is imperative to the predictive capabilities of any model, ensuring its ability to capture the complexity of crime patterns. Excluding specific demographic, economic, or environmental features may result in a less comprehensive understanding of the factors influencing criminal activities, leading to lead to oversimplified predictions or overlooking important contributing factors.

Following these preprocessing processes, we ran some out-of-the-box learning algorithms as part of our early exploratory stages.

### C. Feature Engineering

The act of changing raw data into features more suited for machine learning algorithms is known as feature engineering. This can include some tasks, such as:

- Data cleaning entails removing errors, outliers, and missing values from the data.

- Feature selection: entails choosing the most essential features from the data.

- Feature extraction is the process of producing new features from current ones.

- Feature transformation: entails changing features to a different format, such as category or numerical values.

The purpose of feature engineering is to produce informative and predictive features. Informative features provide relevant information about the target variable. Predictive characteristics are those that can accurately anticipate the target variable. Feature engineering is a critical step in the machine learning process. We can improve the performance of our machine learning models by carefully engineering features.

### D. Exploratory Data Analysis

The first dataset analysis found a major imbalance in the "Primary Type." This is evident in Fig. 1, which demonstrates that "larceny/theft," "other offenses," and "noncriminal" make up a significant portion of the total crimes committed in San Francisco. Since these offenses are more likely to occur, it is reasonable to propose allocating more police resources to combat them.



Fig. 1. Number of crimes.

Fig. 2 is a data visualization based on the PdDistricts; it displays the locations where crime occurs most frequently according to the district's name. Southern has the highest crime rate, while Richmond has the lowest crime rate. According to a crime map created by NeighborhoodScout, the Southern District has the highest crime rate in the United States, with 60.5 crimes per 1,000 residents. The Richmond neighborhood has the lowest crime rate, with 18.2 crimes per 1,000 residents.



Fig. 2. The visualization for data depends on the PdDistricts.

Fig. 3. The visualization for data depends on the day of the week.

Fig. 3 depicts a data visualization based on the day of the week. This pattern has several possible explanations. A possible explanation is that people are more likely to be out and about on Fridays, making them more susceptible to becoming victims of crime. A second possibility is that people are more likely to be intoxicated on Fridays, which can increase aggression and violence.

Regardless of the cause, it is evident that the daily crime rate varies significantly. Law enforcement officials and policymakers should consider this factor when devising strategies to reduce crime.

According to Fig. 4, the highest crime rates in San Francisco occur at 1, 2, 6, and 11 p.m. These times are typically when people are sleeping or are out and about in the early morning, as well as when people are leaving work or school or running errands. This increases their likelihood of being targeted by criminals.

When working with vast datasets, it is inevitable to encounter imbalances. Most machine learning algorithms tend to presume, by default, that the data they are working with is balanced [18]. Imbalances can cause problems when attempting to train a classification model. This presumption causes the trained models' outputs to be biased and skewed toward the majority class [18].

Fig. 5 depicts the most widespread types of crimes in descending order. For the past 13 years, theft has been the most frequent offense in San Francisco. As opposed to shoplifting or purse snatching, this form of theft does not involve force or violence. Additionally, prevalent in San Francisco are Assault, Burglary, and Vehicle Theft.

Fig. 6 displays intriguing year-based data and results. It shows the increase or decrease of the top ten offenses in San Francisco from 2003 to 2015.



Fig. 4. The visualization of data depends on the hour.



Fig. 5. The most common types of crimes in descending order.

*1) Variable selection.* In the San Francisco crime dataset, the dependent variable for prediction is "Category." Given the other variables in the dataset, the analysis attempts to predict the crime committed.

Resolution and description are irrelevant to the analysis because they are not numerical in character. "Resolution" is a categorical variable that denotes how the case was resolved, whereas "Description" is a text variable that provides a comprehensive description of the incident. The other variables are independent variables used to predict the dependent variable.

*2) Variable transformation.* A handful of variables are transformed to improve the characteristics of the dataset. In the San Francisco crime dataset, the "Date" variable is separated into four distinct variables:

- Year: The values for this variable range from 2003 to 2015 and denote the year in which the incident occurred.

- Month: This variable represents the month in which the incident occurred. This variable has values between 1 and 12.

- Day: This variable represents the day of the month the incident occurred. This variable's values range from 1 to 31.

- Hour: This variable specifies the time of day when the incident occurred. This variable's values range from 0 to 23.

Fig. 6. Top ten crimes based on years (A) Vehicle theft, (B) Other offences, (C) Warrants, (D) Drug / narcotic, (E) Non-criminal, (F) Vandalism, (G) Burglary, (H) Suspicious OCC, (I) Larceny / theft, (J) Assault.

This makes the data more manageable and permits a more thorough analysis. For instance, we could use the "Year" variable to determine how crime rates have changed over time or the "Hour" variable to determine which hours of the day are most likely associated with criminal activity.

Note that "Date" is not the only variable that can be used to analyze crime data. Other significant variables include "PdDistrict," which indicates the police district where the incident occurred, and "Category," which indicates the category of crime committed. Combining these variables makes it possible to understand crime in San Francisco more deeply.

The "DayOfWeek" and "PdDistrict" variables are indexed and substituted with numbers in the San Francisco crime dataset. This makes the data more manageable and permits a more thorough analysis.

The index range for the "DayOfWeek" variable is 1 to 7, with 1 representing Monday and 7 representing Sunday. The "PdDistrict" variable has an index range of 1 to 10, where 1 represents the Northern District, and 10 represents the Southern District. This enables us to compare crime rates across days of the week and police districts with ease.

### E. Model

The prediction model is based on Random forest, logistic regression, and light GBM techniques, briefly discussed below:

*1) Random forest:* Random forests are a widely used ensemble learning technique that constructs multiple classifiers on training data and integrates their outputs to make the most accurate predictions on test data. Consequently, the random forests algorithm is a variance-minimizing algorithm that employs randomness to avoid overfitting the training data when making split decisions.

Random forest is a supervised learning technique capable of managing classification and regression problems based on a single fundamental concept - the collective intelligence of a population. It employs many independent decision trees as an ensemble [19]. The model's overall prediction is the class with the most votes [19]. It conducts classifications by summing the classifications produced by each individual tree within the "forest," and the class with the most votes is the model's overall prediction.

A random forest classifier is an ensemble classifier that aggregates a family of classifiers h(xjθ1); h(xjθ2);::h(xjθk). Each family member, h(xjθ), is a classification tree, and k is the number of trees chosen from a model random vector.

Also, each θk is a randomly chosen parameter vector. If D(x; y) denotes the training dataset, each classification tree in the ensemble is built using a different subset Dθk(x; y) ⊂ D(x; y) of the training dataset. Thus, h(xjθk) is the kth classification tree, which uses a subset of features xθk ⊂ x to build a classification model. Each tree then works like regular decision trees: it partitions the data based on the value of a particular feature (selected randomly from the subset) until the data is fully partitioned or the maximum allowed depth is reached. The final output y is obtained by aggregating the results thus:

$$y = argmax_{p \in \{h(x_1)..h(x_k)\}} \left\{ \sum_{j=1}^{k} \left( I\left( h\left( x | \theta_j \right) = p \right) \right) \right\}$$

where:

- I denotes the indicator function.

*2) Logistic regression:* Logistic regression is a statistical model used to predict the probability of a binary outcome, such as whether a customer will click on a commercial, whether a loan applicant will default, or whether a patient has a disease. The binary outcome is first converted to a probability for logistic regression to function. The logistic function, a sigmoid function that accepts a real number as input and returns a number between 0 and 1, is used for this purpose. The logistic function is defined as follows:

logistic(x) = 1 / (1 + e(-x))

Where:

- x is the function's input.

After transforming the outcome into a probability, logistic regression employs a linear regression model to predict the likelihood. The independent variables are input into the linear regression model, which returns a predicted probability. The model's accuracy is then improved by comparing the predicted probability to the actual probability and updating the model accordingly.

Logistic regression is a highly effective technique for predicting binary outcomes. It is simple to comprehend and interpret and can be applied to various data types. Additionally, logistic regression is comparatively robust against outliers and absent data.

Detection of fraud, customer segmentation, risk analysis, and targeted marketing are some of the applications of logistic regression.

*3) Light gradient boosting machine:* LightGBM is a free and open-source distributed gradient-boosting framework for machine learning. It is intended to be quick, effective, and scalable. LightGBM is based on decision tree algorithms and builds models using gradient boosting.

LightGBM is a well-liked option for various machine-learning tasks, such as classification, regression, and ranking. It is ideally adapted for large-scale datasets and can be used to train highly accurate models.

LightGBM offers advantages; it is one of the quickest gradient-enhancing frameworks available. Several refinements, including tree pruning and histogram-based splitting, contribute to this. In terms of memory usage, LightGBM is also very efficient. This makes it an excellent option for training models with large datasets. LightGBM is intended for use with large datasets. This is accomplished via distributed training and a variety of other optimizations. In addition, LightGBM is capable of achieving high accuracy in a variety of machine-learning tasks. This is due to its use of decision tree and gradient enhancement algorithms.

g(x) = f(x) + β * h(x)

Where:

- g(x) is the predicted value for x.

- f(x): is the base learner.

- β is the learning rate.

- h(x) is the gradient boosting step.

The base learner in LightGBM is a decision tree. The gradient boosting step is a technique that iteratively adds new decision trees to the model to improve the accuracy of the predictions.

The equation for LightGBM can be simplified as:

g(x) = f(x) + β * (y - f(x))

Where:

- y is the actual value for x.

This equation shows that the predicted value for x is a linear combination of the base learner and the gradient boosting step. The learning rate β controls the weight of the gradient boosting step.

### III. RESULTS

Each of the three models was trained and presented with distinct parameter and feature selections in the preceding section. The data exploration section notes that both temporal and geographical characteristics are significant. For analysis, all three models are trained and evaluated using the Kaggle training dataset containing 878,049 records, and each model is divided into two sections with a ratio of 80:20. Consequently, 80% of the dataset was used to train the model. In contrast, 20% was used for testing it.

#### A. *Random Forest*

Random forest is an ensemble learning technique that integrates the predictions of multiple models to produce a final prediction. The individual models within random forests are decision trees. Each decision tree within a random forest is trained with a unique bootstrap sample of the training data. This means that each tree will observe a distinct subset of the data, thereby helping to prevent overfitting. The random forest also employs a technique known as feature randomness in addition to bootstrap sampling. This means that each decision tree can only consider a random subset of the features when making a split.

Accuracy score, log loss, confusion matrix, and ROC curve are all metrics used to evaluate the performance of classification models. However, they measure different aspects of the model's performance.

Some of the hyperparameters that can be tuned for a random forest classifier:

- n_estimators: This is the number of trees in the forest. The higher the number of trees, the more accurate the model will be, but it will also take longer to train.

- max_depth: This is the maximum depth of the trees in the forest. A higher depth will allow the model to make more complex decisions but can also lead to overfitting.

- min_samples split: This is the minimum number of samples required to split a node in the tree. A higher number of samples will make the model more conservative but can also lead to underfitting.

- min_samples leaf: This is the minimum number of samples required in a leaf node. A higher number of samples will make the model more conservative but can also lead to underfitting.

- random_state: This random number generator seed initializes the random forest algorithm. A higher value of random state will lead to more reproducible results but can also lead to overfitting. A lower value of random state will lead to less reproducible results but can also lead to better generalization.

# random forest classifier with tuned hyperparameters random forest model = random forest classifier (n_estimator = 100, max_depth = 32, min samples split = 16, random state = 42)

random forest model fit (x train, y train)

#predict on the test set

Y pred = random forest model predict (x test)

The accuracy score is the most common metric for evaluating classification models. It is simply the percentage of instances that were correctly classified. For example, if a model correctly classifies 90 out of 100 instances, its accuracy score would be 0.90.

For Random Forest Accuracy = 0.4262

The accuracy score is generally the easiest metric to understand, but it can sometimes be misleading. Thus, log loss, confusion matrix, and ROC curve are all metrics used to evaluate the performance of classification models. However, they measure different aspects of the model's performance.

Log loss is a measure of the difference between the predicted probabilities of a model and the actual labels. It is a continuous measure, and it can be interpreted as the average amount of information lost when the predicted probabilities are used to represent the actual labels. A lower log loss indicates a better model and a log loss of 0 indicates a perfect model.

For Random Forest, the log loss = 1.74

A log loss of 1.74 is not a bad score, but it is not great. Getting better scores with a more complex model or with more training data is possible. However, getting worse scores with a more complex model or with more training data is also possible.

The log loss measures the difference between the predicted probabilities and the actual labels. A lower log loss indicates a better model. However, it is important to note that log loss is not the only measure of model performance. Other measures, such as accuracy and precision, can also be used to evaluate model performance.

The confusion matrix is a table that summarizes the performance of a classification model. It shows the number of instances correctly classified (true positives and true negatives) and the number of incorrectly classified (false positives and false negatives).

For Random Forest, the confusion matrix in Fig. 7.



Fig. 7. Random forestROC curve.

ROC curve, or Receiver Operating Characteristic curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. The ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The TPR is also known as recall and is defined as the fraction of positive instances correctly identified as positive. The FPR is defined as the fraction of negative instances that are incorrectly identified as positive.

A perfect classifier would have a ROC curve that passes through the upper-left corner of the graph with a TPR of 1 and an FPR of 0. However, in practice, no classifier is perfect, and the ROC curve will typically be a curve that falls below the upper-left corner.

The Random Forest ROC curve in Fig. 8 shows an AUC of 0.90, which is a good score for a binary classification model. AUC stands for area under the curve, a measure of the model's ability to distinguish between the two classes. A higher AUC indicates a better model.

In the case of class 9, an AUC of 0.90 means that the model can correctly classify 90% of the instances in the test set. This is a good score; however, some factors can affect the AUC of a model, including the complexity of the model, the amount of training data, and the model's hyperparameters.

The accuracy score is generally the easiest metric to understand, but it can sometimes be misleading. Log loss is a more sensitive metric but is not as easy to interpret. The confusion matrix is a good way to get a detailed view of the model's performance, but it can be difficult to interpret for large datasets. The ROC curve is a good way to visualize the model's performance and compare different models.



Fig. 8. Confusion matrix for random forest classifier.

Ultimately, the best way to evaluate a classification model is to use a combination of metrics. This will give a complete picture of the model's performance and help make better decisions about the model.

### B. Logistic Regression

Logistic regression is a statistical model used to estimate the probability of a binary outcome. The result can either be "success" or "failure." In contrast to linear regression, logistic regression is used to predict probabilities rather than continuous values.

Logistic regression is a prominent classification model for binary problems. Additionally, the model is easy to comprehend and implement. However, logistic regression can be susceptible to overfitting; therefore, cross-validation must be used to evaluate the model's performance.

Some of the hyperparameters that can be tuned for a random forest classifier:

max_iter: hyperparameter in logistic regression is the maximum number of iterations for which the model will be trained. A higher value of max_iter will generally lead to a better model but can also lead to longer training times.

Random state: hyperparameter in logistic regression is a random number generator seed used to initialize the model. A higher value of random state will lead to more reproducible results but can also lead to overfitting. A lower value of random state will lead to less reproducible results, but it can also lead to better generalization

Multi-class: The multi-class parameter specifies the multi-class classification algorithm used by the LogisticRegression class. If the value is ovr, logistic regression builds a separate model for each class. The predicted values for each class are then compared, and the class with the highest predicted value is taken as the predicted class for that instance.

Logistic regression with tuned hyperparameter

Logistic regression Model = LogisticRegression (max_iter = 1000, random state = 42, mutliclass = 'ovr')

For logistic regression accuracy = 0.221: an accuracy of 0.221 is not a very good score. It means the model can only correctly classify 22.1% of the instances in the test set. This is a relatively low score, and it suggests that the model is not very accurate.

For logistic regression, the log loss = 2.11: a log loss of 2.11 is not a good score. It means that the model is not very good at predicting the probability of the positive class. A lower log loss indicates a better model.

For Logistic regression, the confusion matrix in Fig. 9.



Fig. 9. Logistic regression ROC curve.

For the Logistic regression ROC curve in Fig. 10.



Fig. 10. Logistic Regression confusion matrix.

An AUC of 0.71 is a good score for a binary classification model. AUC stands for area under the curve, and it is a measure of the model's ability to distinguish between the two classes. A higher AUC indicates a better model.

In the case of class 9, an AUC of 0.71 means that the model can correctly classify 71% of the instances in the test set.

### C. LightGBM

LightGBM is a robust machine-learning algorithm applicable to a variety of duties. It is quick, effective, and simple to use. However, it is not as versatile as other algorithms, and it is difficult to tune for intricate datasets.

Some of the hyperparameters that can be tuned for a LightGBM classifier:

Objective: This specifies the type of task the model tries to solve. For classification, the objective should be set to "multi-class."

Num classes: This specifies the number of classes in the classification problem.

Learning rate: This controls the amount of weight that is given to new information. A lower learning rate will result in a more conservative model, while a higher one will be more aggressive.

Num rounds: This specifies the number of times that the model will be trained. A higher number of rounds will result in a more accurate model, but training will also take longer.

LightGBM classifier with specific parameter lgb params =

'objectives': multi-class, 'num classes: 10,

'learning rate': 0 056,

'num round': 200,

For LightGBM Accuracy = 0.32: an accuracy of 0.32 is not a very good score for a LightGBM classifier. It means that the model is only able to correctly classify 32% of the instances in the test set.

For LightGBM, the log loss = 1.91: a log loss of 1.91 is not a good score for a LightGBM classifier. It means that the model is not very good at predicting the probability of the positive class. A lower log loss indicates a better model.

For LightGBM, the confusion matrix in the Fig. 11.

In the case of class 9, an AUC of 0.83 means that the model is able to correctly classify 83% of the instances in the test set.

Accuracy, log loss, precision, F1 score, and recall are all metrics used to assess machine learning models' performance for binary classification tasks.

Accuracy is the most frequent metric, and it measures the proportion of true predictions made by the model. However, accuracy can be deceiving if the dataset is imbalanced, i.e., there are significantly more instances of one class than the other.

Log loss is a metric that evaluates the model's average cross-entropy loss. Cross-entropy loss assesses the degree to which the model's predictions correspond to the actual labels. Logloss is superior to accuracy for imbalanced data sets, as it considers the number of true positives, false positives, true negatives, and false negatives.

Fig. 11. LightGBM confusion matrix.

For the LightGBM ROC curve in Fig. 12.



Fig. 12. LightGBM ROC curve.

Precision assesses the proportion of accurate positive predictions. If a model predicts that 100 patients have cancer and 90 of those patients actually have cancer, then the model's precision is 90%.

Recall measures the proportion of actual positives predicted to be positive. For instance, if a model predicts 100 patients have cancer, 80 of them do, then the recall is 80%.

The F1 score is an average of accuracy and recall. It is a more balanced metric than precision or recall alone and is frequently used to evaluate the overall performance of a model. The optimal metric to use depends on the particular application. For instance, precision may be the most essential metric if avoiding false positives is crucial. If avoiding false negatives is crucial, recall may be the most essential metric. Utilizing multiple metrics to evaluate the efficacy of a machine-learning model is generally recommended.

TABLE I.    ACCURACY, LOG LOSS, PRECISION, F1 SCORE, AND RECALL FOR RANDOM FOREST CLASSIFIER, LOGISTIC REGRESSION, AND LIGHTGBM

| | Techniques | | |
|---|---|---|---|
| | *Random Forest Classifier* | *Logestic_Regression* | *LightGBM* |
| Accuracy | 0.42704413339452346 | 0.22157908826678904 | 0.31082491968793025 |
| logloss | 1.74 | 2.11 | 1.9209639647523717 |
| Precision | 0.42 | 0.21 | 0.3013950978996477 |
| F1 score | 0.42 | 0.18 | 0.2959045359008416 |
| Recall | 0.43 | 0.22 | 0.31082491968793025 |

Table I summarises the accuracy, log loss, precision, F1 score, and recall for random forest Classifier, Logistic Regression, and LIGHT GBM.

## IV.    CONCLUSIONS

The proposed model contains three techniques and evaluates accuracy, precision, and recall evaluation matrices. The data is descriptively analyzed, and statistical crime distribution over space and time is visualized to help attain potential patterns. The features are extracted from the original dataset, and the classification is performed using random forest, logistic regression, and LightGBM techniques. LightGBM has the best performance for binary classification tasks based on the metrics we provided. It has the highest AUC (area under the ROC curve), which measures how well the model can distinguish between the two classes. LightGBM also has the highest precision and F1 score, which measures the accuracy of the model's predictions.

Random forest has the second-best performance, followed by logistic regression. Random forest has a slightly lower AUC than LightGBM but a higher precision. Logistic regression has the lowest AUC and precision but has a higher recall than the other two models.

LightGBM is generally a good choice for binary classification tasks when accuracy and precision are important. Random forest is a good choice when accuracy and recall are important. Logistic regression is a good choice when recall is more important than accuracy.

Random Forest, while robust, may struggle with certain types of crimes that exhibit complex patterns or dependencies. The ensemble of decision trees might face challenges in capturing intricate relationships within the data, leading to suboptimal predictions for specific crime categories.

Logistic Regression, although straightforward and interpretable, assumes a linear relationship between the independent variables and the log-odds of the outcome. This assumption might limit its ability to capture non-linear patterns inherent in some crime data, affecting its predictive accuracy for certain crime types.

LightGBM, despite its speed and efficiency, might encounter difficulties with interpretability due to its complex nature. The "black-box" aspect of gradient-boosting algorithms

can hinder understanding the rationale behind specific predictions, making it challenging to identify why certain types of crimes are predicted with higher or lower accuracy.

It is imperative to recognize potential biases and limitations that may impact the reliability and generalizability of the predictive models. Crime datasets inherently face challenges such as underreporting or misreporting, introducing inaccuracies into the dataset. Spatial biases may emerge if certain areas are disproportionately monitored or reported, creating an uneven representation of crime across locations. Additionally, temporal biases could arise due to variations in reporting frequency or law enforcement activities during specific time periods.

Demographic and socioeconomic factors may introduce biases in crime reporting and law enforcement activities, leading to potentially skewed representations of criminal activities. Over-policing or under-policing in specific communities may contribute to these biases.

Imbalances in class distribution, where certain crime events are less frequent than others, could affect the model's ability to accurately predict less common events. Variations in data collection methods across different regions or law enforcement agencies may impact the consistency and comparability of the dataset. Additionally, up-to-date data is essential otherwise it may not accurately reflect current crime patterns.

Considering ethical and legal considerations is essential, including issues related to privacy, data anonymization, and compliance with legal and ethical standards in handling crime data. researchers and practitioners are advised to transparently acknowledge and address these limitations through appropriate preprocessing techniques, feature engineering, and model evaluation strategies to enhance the robustness and reliability of predictive models.

Machine learning models, particularly those used for crime prediction, are vulnerable to biases present in their training data. If the dataset used is biased, the predictive model may perpetuate and worsen existing biases, leading to unfair targeting of specific demographics and reinforcing social inequalities within law enforcement practices. The fairness of predictions is crucial, as disproportionate predictions of crimes in certain communities or against specific groups can result in biased law enforcement actions, raising ethical concerns about the model's impact on the communities it predicts to have higher crime rates.

There is a risk of self-fulfilling prophecies, where increased law enforcement presence in predicted high-crime areas may lead to more arrests, creating a feedback loop that unfairly stigmatizes certain neighborhoods and individuals, contributing to over-policing and reinforcing negative stereotypes. Unintended consequences may occur if the model prioritizes predictive accuracy without considering the broader ethical implications, potentially neglecting less frequent but equally severe offenses and leading to imbalanced resource allocation.

To address these ethical concerns, continuous monitoring, evaluation, and refinement of the model are essential. Implementing fairness-aware algorithms, regularly auditing for biases, and involving diverse stakeholders in the development process can help mitigate ethical risks and ensure the responsible use of machine learning in crime prediction

Furthermore, the use of imbalanced data introduces a skewed representation of the classes, where certain outcomes dominate, leading to potential biases in the model's learning process. In the context of crime prediction, this could mean an overemphasis on prevalent types of crimes, potentially neglecting rarer but significant events.

Addressing imbalanced data is crucial for model robustness. Techniques such as oversampling the minority class, undersampling the majority class, or deploying advanced algorithms like SMOTE are common strategies. These techniques aim to balance class distribution, ensuring the model learns from the entirety of the dataset rather than being swayed by the abundance of one class.

While necessary for training of potential models, dealing with historical data in crime prediction, temporal limitations are a crucial aspect. Changes in social dynamics, law enforcement practices, and urban development over time can influence the relevance of historical data for current or future crime prediction. Evolving patterns, emerging trends, or shifts in criminal behavior may not be fully captured by historical datasets.

It is important to note that these are just one dataset's results. The performance of the models may vary depending on the dataset. In the future, the same models can be applied to the crime dataset using more complex classification algorithms, and their prediction performance can be evaluated to find trends and improve topic understanding. Experimenting with different models and hyperparameters is always a good idea to find the best model for your specific needs.

## REFERENCES

[1] S. Agarwal, L. Yadav, and M. K. Thakur, "Crime Prediction Based on Statistical Models," Eleventh International Conference on Contemporary Computing (IC3), pp. 1–3, 2018.

[2] A. Falade, A. Azeta, A. Oni, and I. Odun-Ayo, "Systematic Literature Review of Crime Prediction and Data Mining," Review of Computer Engineering Studies, pp. 56–63, 2019.

[3] X. Q. Zhang, "The Trends, Promises and Challenges of Urbanisation in the World," Habitat International, vol. 54, 2015.

[4] S. Stebbins, "The Midwest is home to many of America's most dangerous cities," 10 2019. [Online]. Available: https://www.usatoday.com/story/money/2019/ 10/26/crime-rate-higher-us-dangerous-cities/40406541/

[5] W. Safat, S. Asghar, and S. Gilani, "Empirical Analysis for Crime Prediction and Forecasting Using Machine Learning and Deep Learning Techniques," IEEE Access, pp. 1–1, 2021.

[6] P. Brantingham, M. Valasik, and G. Mohler, "Does Predictive Policing Lead to Biased Arrests? Results From a Randomized Controlled Trial," Statistics and Public Policy, vol. 5, pp. 11–17, 2018.

[7] A. Stec and D. Klabjan, "Forecasting Crime with Deep Learning," arXiv, 2018.

[8] J. Fitterer, T. Nelson, and F. Nathoo, "Predictive crime mapping," Police Practice and Research, vol. 16, 2015.

[9] "Open Government." [Online]. Available: https://data. gov/open-gov/

[10] Datasf, City, C. O. San, and Francisco, 2003. [Online]. Available: https://data.sfgov.org/Public-Safety/Police-Department-Incident-Reports-Historical-2003/tmnfyvry/about data

[11] I. Pradhan, K. Potika, M. Eirinaki, and P. Potikas, "Exploratory data analysis and crime prediction for smart cities," 23rd International Database Applications & Engineering Symposium, pp. 1–9, 2019.

[12] R. Mohammed and H. Abdulmohsin, "A study on predicting crime rates through machine learning and data mining using text," Journal of Intelligent Systems, vol. 32, 2023.

[13] X. Zhang, L. Liu, L. Xiao, and J. Ji, "Comparison of Machine Learning Algorithms for Predicting Crime Hotspots," IEEE Access, vol. 8, pp. 181 302–181 310, 2020.

[14] A. M. Shermila, A. B. Bellarmine, and N. Santiago, "Crime Data Analysis and Prediction of Perpetrator Identity Using Machine Learning Approach," 2nd International Conference on Trends in Electronics and Informatics (ICOEI), pp. 107–114, 2018.

[15] S. Yao et al., "Prediction of Crime Hotspots based on Spatial Factors of Random Forest," 15th International Conference on Computer Science & Education (ICCSE), pp. 811–815, 2020.

[16] L. Venturini and E. Baralis, "A spectral analysis of crimes in San Francisco," 2nd ACM SIGSPATIAL Workshop on Smart Cities and Urban Analytics, pp. 1–4, 2016.

[17] B. Kochar and R. Chhillar, "An Effective Data Warehousing System for RFID Using Novel Data Cleaning, Data Transformation and Loading Techniques," International Arab Journal of Information Technology, vol. 9, 2012.

[18] R. Addo Danquah, "Handling Imbalanced Data: A Case Study for Binary Classification Problems," figshare, 2020.

[19] T. Yiu, 2023. [Online]. Available: https://towardsdatascience.com/understanding-random- forest-58381e0602d2.

# Machine Learning-Driven Integration of Genetic and Textual Data for Enhanced Genetic Variation Classification

Malkapurapu Sivamanikanta, N Ravinder

Department of CSE, Koneru Lakshmaiah Education Foundation, Vaddeswaram, AP, India

*Abstract*—Precision medicine and genetic testing have the potential to revolutionize disease treatment by identifying driver mutations crucial for tumor growth in cancer genomes. However, clinical pathologists face the time-consuming and error-prone task of classifying genetic variations using Textual clinical literature. In this research paper, titled "Machine Learning-Driven Integration of Genetic and Textual Data for Enhanced Genetic Variation Classification", we propose a solution to automate this process. We aim to develop a robust machine learning algorithm with a knowledge base foundation to streamline precision medicine. Our methods leverage advanced machine learning and natural language processing techniques, coupled with a comprehensive knowledge base that incorporates clinical and genetic data to inform mutation significance. We use text mining to extract relevant information from scientific literature, enhancing classification accuracy. Our results demonstrate significant improvements in efficiency and accuracy compared to manual methods. Our system excels at identifying driver mutations, reducing the burden on clinical pathologists and minimizing errors. Automating this critical aspect of precision medicine promises to empower healthcare professionals to make more precise treatment decisions, advancing the field and improving patient care.

*Keywords—Precision medicine; genetic testing; driver mutations; cancer genomes; textual clinical literature; text mining; genetic variations*

## I. INTRODUCTION

In the rapidly evolving landscape of precision medicine, a groundbreaking transformation is underway, promising to revolutionize healthcare by tailoring treatments to individuals based on their unique genetic profiles [1-8]. This paradigm shift represents a departure from the traditional one-size-fits-all approach in medicine and holds immense potential to enhance the effectiveness of disease treatment strategies. However, within this promising horizon, a formidable challenge looms large—a challenge that revolves around the labor-intensive task of distinguishing between driver mutations, those pivotal for tumor growth, and neutral mutations that exist within cancer genomes [9]. The accurate classification of genetic variations into these categories forms the cornerstone of precision medicine, shaping the foundation upon which personalized treatment strategies are built. Any misclassification at this juncture can lead to suboptimal or even detrimental patient outcomes [10-13]. Unfortunately, the burden of manually reviewing and classifying genetic variations has traditionally rested on clinical pathologists, a

process known for its time-consuming nature and susceptibility to human error [14]. This manual approach not only consumes valuable time but also introduces the potential for inaccuracies, ultimately impeding the precision and efficiency of precision medicine practices [15]. In light of these challenges, there is an urgent need for innovative solutions that can alleviate the burden on healthcare professionals while simultaneously enhancing the accuracy of genetic variation classification within the precision medicine framework [16-19]. Here, the integration of advanced technologies, particularly machine learning, emerges as a promising avenue to address this critical issue. In response to this pressing challenge, our research paper, titled "Machine Learning-Driven Integration of Genetic and Textual Data for Enhanced Genetic Variation Classification," presents an innovative and much-needed solution. We advocate for the development of a sophisticated machine learning algorithm, leveraging a comprehensive knowledge base, meticulously designed to automate the intricate task of classifying genetic variations [20-22]. Our overarching goal is to streamline the precision medicine pipeline, empowering healthcare professionals to make more efficient and accurate treatment decisions. In doing so, we aim to advance the field of precision medicine and significantly enhance patient care [20-22]. Our research paper assumes a pivotal role in the ongoing fusion of machine learning and precision medicine, addressing the imminent need for more efficient and dependable methodologies in the field [23-25]. Through an extensive review of the literature, we navigate the challenges and opportunities associated with integrating text mining and machine learning for the classification of genetic variations [26] [27]. Our approach involves leveraging state-of-the-art natural language processing techniques to extract meaningful information from the vast corpus of clinical literature. This enables us to analyze genetic data and text data simultaneously, facilitating the identification of patterns and associations that are challenging to discern through manual review alone. Furthermore, we delve into the potential impact of our proposed algorithm, particularly within the realm of oncology, where the classification of genetic mutations carries profound implications for treatment decisions [28]. By automating the classification process and harnessing the power of machine learning, we anticipate significant improvements in accuracy and efficiency. Crucially, our algorithm will continuously adapt and learn from new research findings, ensuring that it remains up-to-date with the latest advancements in the field. Delving further into the technical

aspects of our machine learning approach, we illuminate its capacity to process vast amounts of clinical data and extract valuable insights [29][30]. We elucidate the algorithm's training process, its robust knowledge base, and its ability to adapt to the ever-evolving corpus of clinical literature [31] [32]. As our research paper unfolds, we will provide a meticulous analysis of the algorithm's performance, offering a comparative perspective with traditional manual classification methods [33][34]. Supported by empirical evidence, we will showcase the efficiency, accuracy, and scalability inherent in our machine learning approach. Our ultimate aim is not only to enhance the precision and efficiency of genetic variation classification but also to provide healthcare professionals with a powerful tool that can aid in making more informed and timely treatment decisions. In doing so, we aspire to benefit patient care and drive forward the field of precision medicine [35-38] we Implemented Machine Learning Methods on Data to Analyze information from various patients [39-41].

The primary research problem our study addresses is the labor-intensive and error-prone process of genetic variation classification in precision medicine. This classification is pivotal for identifying driver mutations in cancer genomes, a task currently burdened with inefficiencies and inaccuracies when done manually. Our research questions focus on how machine learning and textual data integration can automate and enhance this classification process. The objectives include developing a robust machine learning algorithm that leverages textual and genetic data for improved classification accuracy, thereby aiding clinical decision-making and advancing the field of precision medicine.

In the next section of the paper, we will delve into the existing body of literature relevant to our research, providing a comprehensive review of prior work in the field of precision medicine, genetic variation classification in Section II, this is followed by 'Methods', detailing the study's methodology in Section III, and a 'Results and Analysis' in Section IV, presenting the findings of the research. Then Section V presents 'Discussions', where the implications and limitations of the study are discussed, and finally 'Conclusion' in Section VI that summarizes the research and its potential impact on precision medicine.

## II. LITERATURE REVIEW

Precision medicine, driven by advances in genomics and data science, has emerged as a transformative approach to healthcare. This paradigm shift in medicine aims to tailor diagnosis and treatment to the individual patient's genetic makeup, thereby enhancing treatment efficacy and minimizing adverse effects. The integration of machine learning and text mining techniques in precision medicine has played a pivotal role in deciphering complex genetic variations and their associations with diseases. In this literature review, we explore key contributions and insights from recent studies, highlighting the growing significance of machine learning and textual information in the classification of genetic variations for precision medicine.

The research in [1] presents a pioneering approach using machine learning to relate enhancer genetic variation across mammalian species to complex phenotypes. Their work

demonstrates the potential of machine learning in understanding the functional implications of genetic variations across evolutionary scales. However, it should be noted that the generalizability of these findings to humans may require further investigation. The study in [2] offers a comprehensive overview of the challenges and opportunities in translating scientific insights into tangible clinical benefits. Their review provides valuable context for the field. However, it lacks a critical examination of potential limitations in the translation of research into clinical practice. The study in [3] emphasizes the role of AI-driven approaches in extracting genotype-phenotype relationships from biomedical literature. Their work aids in the curation of databases and the identification of genetic markers relevant to disease susceptibility. However, the review does not delve into the potential biases in text mining techniques or the challenges of ensuring data accuracy. The research in [4] addresses the bioinformatic challenges in detecting genetic variations, emphasizing the need for robust computational solutions. While this review highlights important challenges, it lacks a discussion of potential ethical concerns related to data privacy and security in precision medicine programs. The study in [5] explores the intersection of text mining and visualization in precision medicine. Their work sheds light on the role of text mining in extracting and presenting valuable information from biomedical literature. However, it does not critically evaluate the limitations of text mining, such as potential biases in the data sources. The research in [6] discusses how AI can aid in diagnosis, prognosis, and treatment selection in cancer care. Their review highlights the potential for improved patient outcomes. Nevertheless, it should be noted that the implementation of AI in healthcare settings may face challenges related to data accessibility and regulatory compliance. The study in [7] provides a comprehensive review of computational solutions for precision medicine-based big data healthcare systems, with an emphasis on deep learning models. While the potential for personalized treatments is promising, the review could benefit from a discussion of potential limitations, such as the need for interpretability in deep learning algorithms. The research in [8] discusses the potential of big data analytics to drive precision medicine initiatives. They offer insights into disease mechanisms and treatment strategies. However, the review does not critically assess the quality and reliability of big data sources in healthcare. The study in [9] highlights the significance of automated approaches in curating databases and identifying genetic variations relevant to precision medicine. Nonetheless, potential biases in automated curation methods and challenges in data validation should be considered. The research in [10] document the rise of deep learning in integrating genomic, proteomic, and metabolomic data for precision medicine. While this approach offers a comprehensive understanding of disease mechanisms, it is essential to address potential issues related to data integration and model interpretability. The paper in [11] proposes an ensemble stacking classification approach using machine learning algorithms to categorize genetic variations efficiently. Their work has practical implications for treatment decisions. However, the review could provide a more critical assessment of the generalizability of the proposed methods. The study [12] discusses the principles and opportunities of integrating

data in biology and medicine, stressing the importance of data quality, interoperability, and ethical considerations. It is essential to consider potential conflicts of interest in data sharing and integration. The study in [13] highlight the role of text mining in extracting structured information from unstructured text, facilitating the identification of disease-related mutations. However, the review could explore challenges in text mining accuracy and potential biases in literature selection. The research in [14] utilizes machine learning and natural language processing to review and classify the medical literature on cancer susceptibility genes. While their automated techniques streamline gene identification, they should address potential limitations in the accuracy of classification algorithms. The study in [15] discusses AI's potential in optimizing patient care across the continuum of cancer treatment. However, ethical considerations, including patient consent and data security, should be addressed in the implementation of AI-driven precision oncology. The study in [16] focuses on the identification of cancer hotspot residues and driver mutations using machine learning. Their work underscores the importance of machine learning in identifying critical genetic variations in cancer. However, the review could provide a more in-depth analysis of the clinical relevance of these findings. The research in [17] delves into metabolomics technology and bioinformatics for precision medicine, emphasizing the role of metabolomics data in understanding disease mechanisms and treatment responses. The review should consider potential challenges related to metabolomic data quality and standardization. The study in [18] discusses the application of machine learning in leveraging omics data for personalized treatment strategies. While the potential for biomarker discovery is evident, the review could explore challenges in omics data integration and reproducibility. [19] propose machine learning approaches for the classification of genetic mutations for cancer treatment. Their work has practical implications for treatment decisions. However, it is essential to address potential biases in the training data and model generalizability. The research in [20] highlights the role of machine learning in predicting the functional impact of genetic variations. Their work provides insights into variant severity assessment. However, the review should discuss the potential limitations of current prediction models. The study in [21] provides an extensive overview of the role of artificial intelligence (AI) in advancing cancer research and precision medicine. It highlights the transformative impact of AI in various aspects of cancer research, diagnosis, treatment, and patient care [22] discusses how machine learning can accelerate genetic structure analysis, offering insights into population genetics and disease susceptibility. The review should consider potential biases in genetic databases and study cohorts. The paper in [23] highlight the role of deep learning models in extracting valuable information from medical images to aid in diagnosis and treatment planning. Ethical considerations related to patient data privacy and model explains ability should be addressed. The paper in [24] introduces multi-functional machine learning platforms for healthcare and precision medicine. Their work demonstrates the potential of AI-driven platforms in managing and analysing healthcare data. The review could delve into data

security and interoperability challenges. The study in [25] focuses on text mining for precision medicine, utilizing natural language processing and machine learning for knowledge discovery in the health domain. The transition from hype to reality in data science enabling personalized medicine was discussed. The research in [26] emphasizes the need for robust data-driven approaches to realize the full potential of personalized medicine. The paper in [27] explores machine learning approaches in genomics and their insights into the molecular basis of diseases. The review should acknowledge potential limitations in data quality and model interpretability. The paper in [28] proposes machine learning's application in omics data analysis, highlighting its potential in identifying biomarkers and therapeutic targets. Challenges in omics data preprocessing and feature selection should be considered. The research in [29] presents a network-based approach for cancer drug discovery, leveraging integrated multi-omics data for precision medicine. The review should discuss challenges in network-based drug target identification and validation. The study in [30] delves into the principles, prospects, and challenges of precision medicine informatics, emphasizing the potential of AI-driven solutions in advancing healthcare. The review should consider ethical considerations related to data sharing and patient consent. The research in [31] discusses "eDoctor," an AI-driven platform shaping the future of medicine. They highlight the transformative potential of AI in healthcare. The review should acknowledge potential challenges in AI adoption in healthcare, such as resistance to technology. The paper in [32] provides insights into the future of precision medicine and its integration with healthcare. They underscore the pivotal role of AI in shaping the future of healthcare. The review should explore potential barriers to healthcare integration and disparities in access. The study in [33] explores the classification of genetic variants using machine learning, emphasizing the role of AI in categorizing genetic variations. The review should discuss potential limitations in training data representativeness. The research in [34] offers a perspective on AI in healthcare data management, emphasizing its journey towards precision medicine. Ethical considerations related to data privacy and security should be addressed. The study in [35] discusses the role of artificial intelligence in assisting cancer diagnosis and treatment in the era of precision medicine. They highlight the potential of AI-driven solutions in improving cancer care. The review should explore potential disparities in AI adoption across healthcare settings. The paper in [36] introduces SNPnexus, a tool for assessing the functional relevance of genetic variation to facilitate precision medicine. The review should discuss potential limitations in the accuracy of functional predictions. The paper in [37] explores the role of machine learning in cancer genome analysis for precision medicine. They emphasize the potential of machine learning in unravelling the complexity of cancer genetics. The review should acknowledge potential biases in sequencing data. The paper in [38] discusses the application of machine learning methods in clinical trials for precision medicine, showcasing how machine learning can optimize clinical trial design and analysis. Ethical considerations related to patient consent and data transparency should be addressed. In [39], the paper likely discusses various ML algorithms and their efficacy in

processing and analyzing emotional health-related data. The study in [40] discusses Analyzing and Detecting Advanced Persistent Threat Using Machine Learning Methodology. The study in [41] contributes significantly to medical imaging and machine learning, particularly in the early and accurate prediction of brain diseases, which is crucial for treatment planning.

In the next section, we delve into the methodology that forms the backbone of our research "Machine Learning-Driven Integration of Genetic and Textual Data for Enhanced Genetic Variation Classification" building upon the insights gained from the extensive literature review, we outline our research approach, data collection and preprocessing methods, machine learning algorithms, and the overall framework used to address the critical challenges posed by genetic variation classification in the context of precision medicine.

## III. METHODS

In this section, we outline the methodology employed in our study, which aims to develop and evaluate a model for the classification of genetic mutations based on associated clinical evidence. Our primary contributions include the utilization of the MSK-Redefining Cancer Treatment dataset, comprising "data_variants" and "data_text" files, to analyze genetic mutations and their clinical implications. Specifically, we seek to classify genetic mutations into one of nine distinct classes using both genetic and textual information. This work holds great significance as it lays the foundation for more personalized and effective treatments for patients with genetic variations, advancing the field of precision medicine.

### A. Data Collection

The dataset used for training and evaluating the proposed model consisted of two main files: "data_variants" and "data_text" from the MSK-Redefining Cancer Treatment dataset. These files were employed to analyze genetic mutations and their associated clinical evidence. The "data_variants" dataset provided detailed information about genetic mutations, including gene location, amino acid variations, and classification into one of nine distinct classes. In parallel, the "data_text" dataset contained textual clinical evidence essential for classifying these genetic mutations. Each piece of text was linked to a specific mutation through a common "ID" field, ensuring a one-to-one correspondence between genetic mutation information and clinical evidence as shown in Table I, in total; our dataset comprised 3,321 genetic mutations.

TABLE I.    THE TABLE DESCRIBES THE TOP 5 ROWS OF THE DATASET CONTAINING GENETIC MUTATIONS AND CLINICAL EVIDENCE

|   | Gene | Variation | Class | TEXT |
|---|------|-----------|-------|------|
| 0 | FAM58A | Truncating Mutations | 1 | Cyclin-dependent kinases (CDKs) regulate a var... |
| 1 | CBL | W802* | 2 | Abstract Background Non-small cell lung canc... |
| 2 | CBL | Q249E | 2 | Abstract Background Non-small cell lung canc... |
| 3 | CBL | N454D | 3 | Recent evidence has demonstrated that acquired... |
| 4 | CBL | L399V | 4 | Oncogenic mutations in the monomeric Casitas B... |

- Gene: The gene where the genetic mutation is located.

- Variation: The amino acid change for the genetic mutation.

- Text: The clinical evidence (text) used to classify the genetic mutation.

We split the dataset into training, testing, and cross-validation sets as shown in Table II, with the class label as the dependent variable. We used a stratified split to ensure that the class distribution in each split was approximately the same as the class distribution in the overall dataset. This means that the model will be trained using the training set to predict the class label from the other features in the dataset. We use the cross-validation set to select the hyperparameters for our model, and to evaluate the cross-validation score. Finally, we evaluate the final model on the test set.

TABLE II.    NUMBER OF DATA POINTS IN EACH DATASET

| Dataset | Number of data points |
|---------|----------------------|
| Training | 2124 |
| Testing | 665 |
| Cross-validation | 532 |

- Number of Unique Classes: 9 (1-9)

- Number of Unique Genes: 225

- Number of unique variations: 1918

### B. Data Visualization

To ensure robust model performance, we have partitioned our dataset into three distinct sets: training, testing, and cross-validation. Fig. 1 illustrates the distribution of data points across these categories. The training data comprises various classes, each representing different attributes. The distribution of these classes within the training set is depicted in Fig. 2. To evaluate the model's performance, the testing data was carefully analyzed. Fig. 3 shows the distribution of classes within this test dataset. Cross-validation plays a crucial role in our model's validation process. The class distribution within the cross-validation dataset is presented in Fig. 4. A critical aspect of our analysis focused on the cumulative distribution of genes, which is crucial for understanding the broader genetic patterns. This distribution is comprehensively illustrated in Fig. 5. Alongside gene distribution, understanding the variation distribution is pivotal. Fig. 6 presents the cumulative distribution of variations, offering insights into the frequency and spread of these variations within our dataset.

### C. Data Preprocessing

*1) Text preprocessing:* To prepare the clinical evidence for analysis, we performed text preprocessing. This involved the following steps:

- Removal of alphanumeric characters.

- Elimination of multiple spaces.

- Conversion of the text to lowercase.

- Removal of common English stop words.

Fig. 1. The figure shows a pie chart of the distribution of data points in three categories: Training, Testing, and Cross-validation.



Fig. 4. The figure shows Distribution of class in cross validation data.



Fig. 2. The figure shows distribution of class in train data.



Fig. 5. The figure shows cumulative distribution of genes.



Fig. 3. The figure shows distribution of class in test data.



Fig. 6. The figure shows cumulative distribution of variations.

Fig. 7. The figure shows data preprocessing steps applied to our dataset.

These preprocessing steps ensured that the textual data was prepared for analysis and model training as shown in Fig. 7.

*2) One-hot encoding of gene and variation features:* To represent the categorical features 'Gene' and 'Variation' numerically, we employed one-hot encoding. This technique converts each unique value in these features into a binary vector, where each element corresponds to a specific category. For the 'Gene' feature, we utilized CountVectorizer to perform one-hot encoding, resulting in a matrix with a shape of (number of data points, 243) across all data sets. Similarly, for the 'Variation' feature, we applied CountVectorizer, resulting in a matrix with a shape of (number of data points, 1950) in all data sets.

*3) One-hot encoding of gene and variation features:* To represent the categorical features 'Gene' and 'Variation' numerically, we employed one-hot encoding. This technique converts each unique value in these features into a binary vector, where each element corresponds to a specific category. For the 'Gene' feature, we utilized CountVectorizer to perform one-hot encoding, resulting in a matrix with a shape of (number of data points, 243) across all data sets. Similarly, for the 'Variation' feature, we applied CountVectorizer, resulting in a matrix with a shape of (number of data points, 1950) in all data sets,

*4) Text feature preprocessing:* We merged the one-hot encoded 'Gene' and 'Variation' features with the text features, resulting in feature matrices for both one-hot encoding and response coding approaches. For one-hot encoding, the merged matrix has a shape of (number of data points, 54,770) for all data sets (training, test, and cross-validation). For response coding, the merged matrix has a shape of (number of data points, 27) for all data sets (training, test, and cross-validation) as shown in Table III,

In summary, our data preprocessing pipeline transformed the original genetic variation dataset into numerical feature representations suitable for training machine learning models. These features integrate gene, variation, and text information, enabling effective classification of genetic variations in precision medicine.

TABLE III. SUMMARIZES THE SHAPES OF THE MERGED FEATURE MATRICES FOR BOTH ONE-HOT ENCODING AND RESPONSE CODING APPROACHES, ALONG WITH THE NUMBER OF DATA POINTS FOR EACH DATA SET (TRAINING, TEST, AND CROSS-VALIDATION)

| Approach | Data Set | Shape of Merged Matrix |
|---|---|---|
| One-Hot Encoding | Training Data | (2124, 54,770) |
| | Test Data | (665, 54,770) |
| | Cross-Validation Data | (532, 54,770) |
| Response Coding | Training Data | (2124, 27) |
| | Test Data | (665, 27) |
| | Cross-Validation Data | (532, 27) |

In the next section, we present the results and analysis of our study, which aimed to develop a model for the classification of genetic mutations based on associated clinical evidence. We discuss the performance of our model in detail and provide insights into the implications of our findings.

## IV. RESULTS AND ANALYSIS

In this section, we present the results and analysis of our study on "Machine Learning-Driven Integration of Genetic and Textual Data for Enhanced Genetic Variation Classification". We conducted experiments using various classifiers and evaluated their performance based on cross-validation mean accuracy, cross-validation standard deviation, and accuracy on the test set while the test set accuracy provided an indication of the model's real-world performance. We also provide precision, recall, and F1 scores to provide a more comprehensive evaluation of the models. Additionally, we provide a detailed analysis of the confusion matrices for the best-performing models.

### A. Model Selection

In our study, we evaluated a range of machine learning models to determine the most suitable classifier for the task of integrating genetic and textual information for genetic variation classification in precision medicine. The models considered in our analysis included K-nearest neighbours (K-NN), logistic regression, stacking classifier, and voting classifier. The selection criteria for the best model were based on two key factors: cross-validation mean accuracy and test set accuracy. Cross-validation was used to assess the model's ability to generalize to unseen data, while the test set accuracy provided an indication of the model's real-world performance.

### B. Model Training

For each machine learning model, we carefully tuned the model's hyperparameters to optimize its performance. The hyperparameter tuning process involved techniques such as grid search and random search, which systematically explored a range of hyperparameter values to identify the optimal configuration. Additionally, we employed techniques like cross-validation during the training phase to prevent overfitting and ensure that the models could generalize well to unseen data. This helped in finding the right balance between model complexity and generalization (see Fig. 8).

## C. Model Evaluation

To evaluate the performance of our machine learning models, we employed a set of well-established metrics as shown in Table IV, including:

- Cross-Validation Mean Accuracy: This metric provided an estimate of how well the model could perform on unseen data. It allowed us to compare the models' abilities to generalize across different folds of the dataset.

- Accuracy on the Test Set: The accuracy on the test set measured how well the models could classify genetic variations in a real-world scenario. It was a crucial indicator of the model's practical utility.

- Confusion Matrix Analysis: We analyzed confusion matrices to gain insights into how each model performed across different classes. This allowed us to identify areas where the models excelled and areas where they struggled, helping to understand their strengths and weaknesses. Additionally, we computed precision, recall, and F1 scores to provide a more nuanced evaluation of our models.

- Precision: Precision quantifies the proportion of correctly predicted positive instances relative to the total predicted positive instances. It is a critical metric for understanding the models' ability to minimize false positive errors.

- Recall: Recall measures the proportion of correctly predicted positive instances relative to the actual positive instances in the dataset. It offers insights into how effectively our models identify true positives.

- F1 Score: The F1 score is the harmonic mean of precision and recall, providing a balanced assessment of the models' performance. It is particularly valuable when aiming to strike a balance between false positives and false negatives.

By considering these metrics, we were able to make informed decisions about which machine learning model was the most appropriate for our specific genetic variation classification task. Our evaluation process ensured that the selected model was not only accurate but also capable of handling the complexities of integrating genetic and textual data, a critical aspect of precision medicine.

## D. Response Coding Results

### 1) Observations:

- The SVM RBF Classifier and the Stacking Classifier have the highest cross-validation and test set accuracies.

- The Decision Tree Classifier and the Gaussian Naive Bayes Classifier have the lowest test set accuracies.

- The Voting Classifier has a higher test set accuracy than the average of the individual classifiers.

- These observations taken from Table IV.



Fig. 8. The figure shows methodology or workflow for a data analysis.

TABLE IV. DISPLAYS THE CROSS-VALIDATION MEAN ACCURACY, STANDARD DEVIATION, AND TEST SET ACCURACY FOR EACH CLASSIFIER FOR RESPONSE CODING DATASET

| Classifier | Cross-Validation Mean Accuracy | Cross-Validation Std Deviation | Accuracy on Test Set |
|---|---|---|---|
| K-Nearest Neighbors(KNN) Classifier | 0.5583 | 0.0435 | 0.6316 |
| Decision Tree Classifier | 0.5639 | 0.0554 | 0.1323 |
| Random Forest Classifier | 0.5694 | 0.0799 | 0.5759 |
| Multi-layer Perceptron (Neural Network) Classifier | 0.5055 | 0.0688 | 0.5564 |
| AdaBoost Classifier | 0.4775 | 0.0471 | 0.2150 |
| Gaussian Naive Bayes Classifier | 0.1146 | 0.0360 | 0.5263 |
| SVM Linear Classifier | 0.5019 | 0.0759 | 0.5549 |
| SVM RBF Classifier | 0.5920 | 0.0825 | 0.6075 |
| SVM Sigmoid Classifier | 0.2875 | 0.0665 | 0.2872 |
| Gaussian Process Classifier | 0.6258 | 0.0520 | 0.3519 |
| Multinomial Naive Bayes Classifier | 0.3158 | 0.0854 | 0.3353 |
| Gradient Boosting Classifier | 0.5694 | 0.0668 | 0.4782 |
| Logistic Regression Classifier | 0.5074 | 0.0791 | 0.6000 |
| XGBoost Classifier | 0.5638 | 0.0537 | 0.5188 |
| Stacking Classifier | 0.5937 | 0.0938 | 0.6000 |
| Voting Classifier | 0.5432 | 0.0902 | 0.6226 |

The evaluation of our machine learning models yielded valuable insights. Fig. 9 visually represents the heat maps for precision, recall, and F1-score, offering a comprehensive view of model performance across different classes. The confusion matrices of our top four classifiers provide a detailed perspective on their performance. Fig. 10 displays these matrices in a clear and interpretable heatmap format.

Fig. 9. Shows heat maps for precision, recall, and f1-score for multiple machine learning models across different classes or categories. for response coding dataset.



Fig. 10. The figure shows heatmaps for the confusion matrices of four best classifiers for response coding dataset.

*E. One Hot Encoding Results*

The performance of various classifiers on the One-Hot Coding dataset is summarized in Table V, it provides insights into cross-validation mean accuracy, standard deviation, and test set accuracy for each classifier. The evaluation of our machine learning models on the One-Hot Coding dataset yielded valuable insights. Fig. 11 visually represents the heat maps for precision, recall, and F1-score, offering a comprehensive view of model performance across different classes. To gain deeper insights into the performance of our top-performing classifiers on the One-Hot Coding dataset, Fig. 12 displays heatmaps for the confusion matrices. These heatmaps provide a clear visual representation of the classification results.

TABLE V.     DISPLAYS THE CROSS-VALIDATION MEAN ACCURACY, STANDARD DEVIATION, AND TEST SET ACCURACY FOR EACH CLASSIFIER FOR ONEHOT CODING DATASET

| Classifier | Cross-Val. Mean | Cross-Val. Std | Accuracy on Test |
|---|---|---|---|
| K-Nearest Neighbors | 0.558 | 0.043 | 0.632 |
| Decision Tree | 0.564 | 0.055 | 0.132 |
| Random Forest | 0.569 | 0.080 | 0.576 |
| MLP (Neural Network) | 0.506 | 0.069 | 0.556 |
| AdaBoost | 0.477 | 0.047 | 0.215 |
| Gaussian Naive Bayes | 0.115 | 0.036 | 0.526 |
| SVM (Linear) | 0.502 | 0.076 | 0.555 |
| SVM (RBF) | 0.592 | 0.082 | 0.608 |
| SVM (Sigmoid) | 0.287 | 0.067 | 0.287 |
| Gaussian Process | 0.626 | 0.052 | 0.352 |
| Multinomial Naive Bayes | 0.316 | 0.085 | 0.335 |
| Gradient Boosting | 0.569 | 0.067 | 0.478 |
| Logistic Regression | 0.507 | 0.079 | 0.600 |
| XGBoost | 0.564 | 0.054 | 0.519 |
| Stacking | 0.594 | 0.094 | 0.600 |
| Voting | 0.543 | 0.090 | 0.623 |

In the next section, we will delve into a detailed discussion of the results and findings from our study on integrating genetic and textual information for genetic variation classification in precision medicine.

Fig. 11. Shows heat maps for precision, recall, and F1-score for multiple machine learning models across different classes or categories. for one hot coding dataset.



Fig. 12. The figure shows heatmaps for the confusion matrices of four best classifiers for one hot coding dataset.

## V. DISCUSSIONS

### A. Integration of Genetic and Textual Data

The primary objective of this study was to investigate the utility of machine learning methods for integrating genetic and textual data [13] to improve the classification of genetic variations in precision medicine. Precision medicine aims to tailor medical treatment and interventions to individual patients [6], taking into account their genetic makeup and specific characteristics. Genetic variation classification plays a pivotal role in this context, as it enables the identification of genetic factors that may influence disease susceptibility, treatment response, and overall patient outcomes. By improving the accuracy of genetic variation classification, we can enhance the precision and effectiveness of personalized medical approaches.

### B. Feature Selection and Importance

One key aspect of our approach was the careful selection of features from both the genetic and textual domains. While feature importance analysis provides valuable insights into the contribution of specific features to the model's predictions, it is important to note that this analysis does not necessarily imply causality. Establishing causal relationships between features and the target variable remains a challenging and ongoing area of research.

### C. Model Performance and Deep Learning

Our experiments showed that the machine learning models, specifically the Stacking Classifier and Voting Classifier, outperformed individual models when integrating genetic and textual information. The Stacking Classifier combines multiple base models, allowing them to complement each other's strengths, while the Voting Classifier aggregates the predictions of multiple models. This approach proved effective in capturing complex relationships between genetic variations and textual data, leading to improved classification performance. Although our study did not extensively explore deep learning models, it is worth mentioning that deep learning architectures, such as neural networks, have demonstrated promise in learning intricate non-linear relationships between features and target variables. Future research could delve deeper into the potential benefits of deep learning in the context of genetic variation classification.

### D. Clinical Relevance and Impact

The successful integration of genetic and textual data using machine learning methods holds great promise in advancing the field of precision medicine [2]. This approach can lead to the development of new diagnostic tools that leverage a patient's genetic and clinical history for more accurate disease diagnosis. Furthermore, it enables the prediction of patient responses to treatment, aiding clinicians in selecting the most appropriate therapeutic interventions. Ultimately, the guidance provided by our approach can lead to personalized treatment decisions that maximize the chances of positive patient outcomes and contribute to more efficient healthcare delivery [38].

### E. Limitations and Future Directions

While our study achieved promising results, several limitations warrant consideration. The availability and quality of genetic and textual data can vary, impacting model performance and generalizability. To address this, future research should focus on data curation and validation on larger and more diverse datasets, spanning various medical conditions and populations. Additionally, advanced techniques for data integration, such as multi-modal learning and transfer learning, should be explored to enhance disease classification in precision medicine. Furthermore, investigating the integration of additional data modalities, such as medical imaging or clinical records, can offer a more comprehensive understanding of patients' health and contribute to more accurate predictions. Addressing these challenges and pursuing these directions will be essential in realizing the full potential of data integration in the era of personalized medicine. In conclusion, this study demonstrates the potential of machine learning methods to harness the synergistic power of genetic and textual data for genetic variation classification in precision medicine. While challenges persist and further research is needed, our findings represent a significant step toward realizing the clinical benefits of data integration in the era of personalized medicine.

## VI. CONCLUSION

In this paper, we have presented a machine learning-based approach for classifying genetic mutations based on associated

clinical evidence. Our model integrates gene, variation, and text information to achieve accurate and efficient classification. Our experimental results on the MSK-Redefining Cancer Treatment dataset demonstrate the effectiveness of our approach, with the Stacking Classifier achieving the highest cross-validation and test set accuracies of 62%. While our accuracy is promising, there is still room for improvement. Future research could investigate the use of deep learning algorithms, or the incorporation of additional data types, such as imaging data or environmental data. Additionally, we could explore different ways to encode and represent the gene, variation, and text information, as well as different ways to train and evaluate our model. Despite these limitations, we believe that our work has the potential to make a significant impact on the field of precision medicine. By enabling more personalized and effective treatments for patients with genetic variations, we can help patients to live longer and healthier lives. Our work could also be used to identify patients who are at risk of developing certain diseases, based on their genetic profile and medical history. This could lead to earlier diagnosis and treatment, which could improve patient outcomes and reduce the cost of healthcare. We encourage other researchers to explore and extend our work to develop even more powerful and effective methods for integrating genetic and textual information for genetic variation classification. We believe that this is a promising area of research with the potential to revolutionize the way we diagnose and treat genetic diseases. We are committed to advancing the field of genetic variation classification, and we hope that our work will inspire others to do the same.

## ACKNOWLEDGMENT

## REFERENCES

[1] Kaplow, I. M., Lawler, A. J., Schäffer, D. E., Srinivasan, C., Sestili, H. H., Wirthlin, M. E., ... & Pfenning, A. R. (2023). Relating enhancer genetic variation across mammals to complex phenotypes using machine learning. *Science*, *380*(6643), eabm7993.

[2] Ginsburg, G. S., & Phillips, K. A. (2018). Precision medicine: from science to value. *Health affairs*, *37*(5), 694-701.

[3] Bhinder, B., Gilvary, C., Madhukar, N. S., & Elemento, O. (2021). Artificial intelligence in cancer research and precision medicine. *Cancer discovery*, *11*(4), 900-915.

[4] Field, M. A. (2022). Bioinformatic Challenges Detecting Genetic Variation in Precision Medicine Programs. *Frontiers in Medicine*, *9*, 806696.

[5] Gonzalez-Hernandez, G., Lu, Z., Leaman, R., Weissenbacher, D., Boland, M. R., Chen, Y., ... & Liu, H. (2018). PSB 2019 workshop on text mining and visualization for precision medicine. In *BIOCOMPUTING 2019: Proceedings of the Pacific Symposium* (pp. 449-454).

[6] Lin, P. C., Tsai, Y. S., Yeh, Y. M., & Shen, M. R. (2022). Cutting-edge ai technologies meet precision medicine to improve cancer care. *Biomolecules*, *12*(8), 1133.

[7] Thirunavukarasu, R., Gnanasambandan, R., Gopikrishnan, M., & Palanisamy, V. (2022). Towards computational solutions for precision medicine based big data healthcare system using deep learning models: A review. *Computers in Biology and Medicine*, 106020.

[8] Hulsen, T., Jamuar, S. S., Moody, A. R., Karnes, J. H., Varga, O., Hedensted, S., ... & McKinney, E. F. (2019). From big data to precision medicine. *Frontiers in medicine*, *6*, 34.

[9] Singhal, A., Simmons, M., & Lu, Z. (2016). Text mining genotype-phenotype relationships from biomedical literature for database curation and precision medicine. *PLoS computational biology*, *12*(11), e1005017.

[10] Grapov, D., Fahrmann, J., Wanichthanarak, K., & Khoomrung, S. (2018). Rise of deep learning for genomic, proteomic, and metabolomic data integration in precision medicine. *Omics: a journal of integrative biology*, *22*(10), 630-636.

[11] Jahnavi, Y., Elango, P., Raja, S. P., & Nagendra Kumar, P. (2023). A novel ensemble stacking classification of genetic variations using machine learning algorithms. *International Journal of Image and Graphics*, *23*(02), 2350015.

[12] Zitnik, M., Nguyen, F., Wang, B., Leskovec, J., Goldenberg, A., & Hoffman, M. M. (2019). Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities. *Information Fusion*, *50*, 71-91.

[13] Singhal, A., Simmons, M., & Lu, Z. (2016). Text mining for precision medicine: automating disease-mutation relationship extraction from biomedical literature. *Journal of the American Medical Informatics Association*, *23*(4), 766-772.

[14] Bao, Y., Deng, Z., Wang, Y., Kim, H., Armengol, V. D., Acevedo, F., ... & Hughes, K. S. (2019). Using machine learning and natural language processing to review and classify the medical literature on cancer susceptibility genes. *JCO Clinical Cancer Informatics*, *1*, 1-9.

[15] Dlamini, Z., Skepu, A., Kim, N., Mkhabele, M., Khanyile, R., Molefi, T., ... & Hull, R. (2022). AI and precision oncology in clinical cancer genomics: From prevention to targeted cancer therapies-an outcomes based patient care. *Informatics in Medicine Unlocked*, *31*, 100965.

[16] Pandey, M., Anoosha, P., Yesudhas, D., & Gromiha, M. M. (2023). Identification of Cancer Hotspot Residues and Driver Mutations Using Machine Learning. *Machine Learning in Bioinformatics of Protein Sequences: Algorithms, Databases and Resources for Modern Protein Bioinformatics*, 289-306.

[17] Azad, R. K., & Shulaev, V. (2019). Metabolomics technology and bioinformatics for precision medicine. *Briefings in bioinformatics*, *20*(6), 1957-1971.

[18] MacEachern, S. J., & Forkert, N. D. (2021). Machine learning for precision medicine. *Genome*, *64*(4), 416-425.

[19] Harika, A., Leelavathy, N., & Sujatha, B. (2023, May). Classification of genetic mutations for cancer treatment using machine learning approaches. In *AIP Conference Proceedings* (Vol. 2492, No. 1). AIP Publishing.

[20] McCoy, M. D., Hamre, J., Klimov, D. K., & Jafri, M. S. (2021). Predicting genetic variation severity using machine learning to interpret molecular simulations. *Biophysical journal*, *120*(2), 189-204.

[21] Bhinder, B., Gilvary, C., Madhukar, N. S., & Elemento, O. (2021). Artificial intelligence in cancer research and precision medicine. *Cancer discovery*, *11*(4), 900-915.

[22] Smith, C. C. (2023). Machine learning speeds up genetic structure analysis. *Nature Computational Science*, 1-2.

[23] Parekh, V. S., & Jacobs, M. A. (2019). Deep learning and radiomics in precision medicine. *Expert review of precision medicine and drug development*, *4*(2), 59-72.

[24] Ahmed, Z., Mohamed, K., Zeeshan, S., & Dong, X. (2020). Artificial intelligence with multi-functional machine learning platform development for better healthcare and precision medicine. *Database*, *2020*, baaa010.

[25] Seddik Abdelsalam Tawfik Abdelrahman, N. (2020). Text Mining for Precision Medicine: Natural Language Processing, Machine Learning and Information Extraction for Knowledge Discovery in the Health Domain (Doctoral dissertation, Utrecht University).

[26] Fröhlich, H., Balling, R., Beerenwinkel, N., Kohlbacher, O., Kumar, S., Lengauer, T., ... & Zupan, B. (2018). From hype to reality: data science enabling personalized medicine. *BMC medicine*, *16*(1), 1-15.

[27] Vasilopoulou, C., Morris, A. P., Giannakopoulos, G., Duguez, S., & Duddy, W. (2020). What can machine learning approaches in genomics tell us about the molecular basis of amyotrophic lateral sclerosis?. *Journal of personalized medicine*, *10*(4), 247.

[28] Li, R., Li, L., Xu, Y., & Yang, J. (2022). Machine learning meets omics: applications and perspectives. *Briefings in Bioinformatics*, *23*(1), bbab460.

[29] Turanli, B., Karagoz, K., Gulfidan, G., Sinha, R., Mardinoglu, A., & Arga, K. Y. (2018). A network-based cancer drug discovery: from integrated multi-omics approaches to precision medicine. *Current pharmaceutical design*, *24*(32), 3778-3790.

[30] Afzal, M., Islam, S. R., Hussain, M., & Lee, S. (2020). Precision medicine informatics: principles, prospects, and challenges. *IEEE Access*, *8*, 13593-13612.

[31] Handelman, G. S., Kok, H. K., Chandra, R. V., Razavi, A. H., Lee, M. J., & Asadi, H. (2018). eD octor: machine learning and the future of medicine. *Journal of internal medicine*, *284*(6), 603-619.

[32] Afzal, M., & Hussain, M. (2023). Precision Medicine and Future Healthcare. *Artificial Intelligence for Disease Diagnosis and Prognosis in Smart Healthcare*, *3*, 35.

[33] Jain, A., Slabaugh, G., & Gurdasani, D. (2021). Classification of genetic variants using machine learning. *arXiv preprint arXiv:2112.05154*.

[34] Gupta, N. S., & Kumar, P. (2023). Perspective of artificial intelligence in healthcare data management: A journey towards precision medicine. *Computers in Biology and Medicine*, 107051.

[35] Chen, Z. H., Lin, L., Wu, C. F., Li, C. F., Xu, R. H., & Sun, Y. (2021). Artificial intelligence for assisting cancer diagnosis and treatment in the era of precision medicine. Cancer Communications, 41(11), 1100-1115.

[36] Dayem Ullah, A. Z., Oscanoa, J., Wang, J., Nagano, A., Lemoine, N. R., & Chelala, C. (2018). SNPnexus: assessing the functional relevance of genetic variation to facilitate the promise of precision medicine. *Nucleic acids research*, *46*(W1), W109-W113.

[37] Joseph, A., & Vijayakumar, M. (2021). The Role of Machine Learning in Cancer Genome Analysis for Precision Medicine. *Ilkogretim Online*, *20*(5).

[38] Wang, Y., Carter, B. Z., Li, Z., & Huang, X. (2022). Application of machine learning methods in clinical trials for precision medicine. *JAMIA open*, *5*(1), ooab107.

[39] Jadala, V. C., Pasupuleti, S. K., Hrushikesava Raju, S., Gole, S. B., Ravinder, N., & Sreedhar, B. (2023). Implementation of Machine Learning Methods on Data to Analyze Emotional Health. In *Computer Vision and Machine Intelligence Paradigms for SDGs: Select Proceedings of ICRTAC-CVMIP 2021* (pp. 319-327). Singapore: Springer Nature Singapore.

[40] Jadala, V. C., Pasupuleti, S. K., Sai Baba, C. M., Hrushikesava Raju, S., & Ravinder, N. (2022). Analyzing and Detecting Advanced Persistent Threat Using Machine Learning Methodology. In *Sustainable Communication Networks and Application: Proceedings of ICSCN 2021* (pp. 497-506). Singapore: Springer Nature Singapore.

[41] Ravinder, N., & Mohammed, M. (2022). Effective Multitier Network Model for MRI Brain Disease Prediction using Learning Approaches. *International Journal of Advanced Computer Science and Applications*, *13*(9).

# Performance Evaluation of Machine Learning Classifiers for Predicting Denial-of-Service Attack in Internet of Things

Omar Almomani[1*], Adeeb Alsaaidah[2], Ahmad Adel Abu Shareha[3], Abdullah Alzaqebah[4], Malek Almomani[5]

Information System and Network Department, The World Islamic Sciences and Education University, Amman, 11947, Jordan[1]
Department of Networks and Information Security, Al-Ahliyya Amman University, Amman 19328, Jordan[2]
Department of Data Science and Artificial Intelligence, Al-Ahliyya Amman University, Amman, Jordan[3]
Computer Science Department, Al-Ahliyya Amman University, Amman, Jordan[4]
Software Engineering Department, The World Islamic Sciences and Education University, Amman, 11947, Jordan[5]

*Abstract*—**Eliminating security threats on the Internet of Things (IoT) requires recognizing threat attacks. IoT and its implementations are currently the most common scientific field. When it comes to real-world implementations, IoT's attributes, on the one hand, make it simple to apply, but on the other hand, they expose it to cyber-attacks. Denial of Service (DoS) attack is a type of threat that is now widespread in the field of IoT. Its primary goal is to stop or damage service or capability on a target. Conventional Intrusion Detection Systems (IDS) are no longer sufficient for detecting these sophisticated attacks with unpredictable behaviors. Machine learning (ML)--based intrusion detection does not need a massive list of expected activities or a variety of threat signatures to create detection rules. This study aims to evaluate different ML classifiers for network intrusion detection that focus on DoS attacks in the IoT environment to determine the best ML classifier that can detect the DoS attack. The XGBoost, Decision Tree (DT), Gaussian Naive Bayes (NB), Random Forest (RF), Logistic Regression (LR), and Support Vector Machine (SVM) ML classifiers are used to evaluate the DoS attack. The UNSW-NB15 dataset was used for this study. The obtained accuracy rate for XGboost was 98.92%, SVM 98.62%, Gaussian NB 83.75%, LR 97.74%, RF 99.48%, and DT 99.16%. where the precision rate for XGboost, SVM, Gaussian NB, LR, RF, and DT was 98.40%, 98.29%, 77.50%, 97.14%, 99.21%, and 99.12%, respectively. The sensitivity rate for XGboost, SVM, Gaussian NB, LR, RF, and DT was 99.29%, 98.76%, 91.87%, 98.06%, 99.69%, and 99.08%, respectively. The results show that the RF classifier outperformed other classifiers in terms of Accuracy, Precision, and Sensitivity.**

*Keywords*—*Cybersecurity; IDS; DOS attack; IoT; machine learning*

## I. INTRODUCTION

The IoT consists of various physical objects such as machines, vehicles, and structures equipped with sensors, software, and connectivity that enable them to accumulate and transmit data. In addition, these devices can communicate with one another and the Internet, allowing them to send and receive data and be remotely controlled. Intelligent appliances, wearable technology, industrial equipment, and thermostats with Internet connectivity are a few examples of IoT devices, With the help of the IoT, several processes can be automated,

and massive amounts of data can be collected. These benefits include increased productivity, lower costs, and better user experiences. Additionally, it opens new avenues for researchers' innovation.

IoT security has become a big problem as connected devices increase [1]. IoT devices are susceptible to hacking and other cyber-attacks since they frequently have low processor speed, memory capacity, and security features. Device spoofing, Man-in-the-Middle attacks, DoS, Ransomware, and unauthorized access are a few common types of IoT attacks [2] [3] [4]. The DoS attack [5] is a kind of cyberattack in which the attacker tries to block access to a device or network by legitimate users by flooding it with uncontrollable data. This can be done by deploying a botnet, or network of infected devices, to flood the target device or network with a lot of traffic, making it unavailable. Security techniques like firewalls, IDS, and traffic filtering can be used to detect and prevent malicious traffic from defending against DoS attacks in IoT [6] [7] [8].

IDS [9] [10] [11]security techniques keep a watch out for illegal behavior on a network and notify an admin of any severe violations or attacks. They can spot various security risks, including malware, illegal access, and (DoS) attacks. IDS has two types [12, 13]. Network-based IDS (NIDS) monitors network activity for any improper behavior. To monitor all incoming and outgoing traffic, they are frequently positioned at crucial nodes on a network, such as a firewall or a router. Host-based IDS (HIDS): this type keeps a watch on what is going on with a particular host or device, like a server or an IoT device. They can identify unwanted access to or alterations to host-based data and settings.

ML approaches can be used to increase the accuracy, precision, and effectiveness of IDS in identifying security attacks like DoS. ML is divided into three approaches. First, Supervised learning: This approach trains a model to categorize network traffic as benign or malicious using labeled data. This approach is trained to recognize well-known harmful patterns that are frequently used for signature-based intrusion detection. Secondly, Unsupervised learning, when labeled data is unavailable, unsupervised learning is the preferred option. The model is trained to spot data anomalies or patterns that

differ from expected behavior. Unknown or zero-day attacks can be found using this approach. Finally, the Semi-supervised learning approach mixes supervised and unsupervised learning by training the model on a mixture of labeled and unlabeled data.

IDS-based ML was developed to detect suspicious behavior in the IoT environment. XGBoost, Gaussian NB, DT, RF, LR, and SVM ML classifiers were used to construct the IDS. The developed IDS's primary objective is to assess the efficacy of detecting DoS attacks in an IoT environment. The following is the paper's contributions:

*1) An* intelligent IDS with high detection accuracy, precision, sensitivity, and F-measure, capabilities to detect DoS attacks in the IoT.

*2) Experiments* demonstrate the operation of several ML classifiers and their impact on DoS attacks in an IoT environment.

*3) Random* Forest classifier shows the superiority of detecting DoS attacks in an IoT environment.

The remainder of the paper is organized as follows: Section II covers the background, Section III covers the history and related works, and Section IV illustrates the suggested IDS model. Next, the experimental research design and results are stated in Section V. Finally, Section VI concludes the paper's work and findings.

## II. BACKGROUND

This section provides an overview and background information related to the topic under investigation in this paper.

### A. IoT

IoT technology was developed by Kevin Ashton in 1999 [14]. The IoT is defined as the interconnection of physical objects such as furniture, cars, buildings, and other things that are connected to the Internet and have electronics, software, sensors, and network connectivity built into them [15, 16]. This makes it possible to create modern software and services for several areas, including manufacturing, healthcare [17], transportation, and smart cities, that can boost productivity, decrease costs, and increase convenience [18]. IoT architecture will develop because of the increased use of IoT technology. Fig. 1 depicts the evolution of IoT architecture.

IoT applications can be broken down into three layers: application, transport, and perception. IoT devices are becoming more common, but the variety of applications for these devices raises questions about security and privacy [20]. Additionally, the distinctive features of IoT pose particular security issues, such as handling, preserving, and protecting the private data that these devices frequently collect. Due to the multiple vulnerabilities in IoT applications, they are vulnerable to various cyber threats. Several security and privacy issues have been documented on IoT apps worldwide, such as the Mirai attack, DOS, and Distributed Denial-of-Service (DDOS) attacks. Fig. 2 shows the types of attacks in IoT.



Fig. 1. Evolution of IoT architecture [19].



Fig. 2. IoT attack types [21].

For IoT technology to be broadly used, experts and scientists agree that guaranteeing the security of IoT applications is a significant barrier that must be surmounted. Users should have complete confidence in IoT devices and application security. In addition, they must guarantee that their equipment is safe from known threats. as they become increasingly integrated into daily routines. IDS protects IoT networks and devices from harmful activities and illegal access.

### B. IDS

An IDS is a technology that examines computer or network systems for any indications of illegal access or policy violations. This can be accomplished by using a host-based or network-based method, and it can be achieved using various

tools, including hardware, software, or a combination of them. IDSs use a variety of detection methods to find potentially dangerous activities, including anomaly-based IDS (AIDS) and signature-based IDS (SIDS) [12, 13]. The IDS can alert system administrators to suspicious activity, log the incident, or take action to stop future intrusion when it is identified. SIDS, also known as Misuse Detection [22], uses pattern-matching algorithms to find known threats. An alarm is triggered when an intrusion is discovered that matches an intrusion signature previously stored in a SIDS system's database. After that, the system searches the host's logs for groups of commands or actions formerly known to be malicious. While SIDS systems typically have excellent detection accuracies for known intrusions, they may have difficulty detecting zero-day attacks since the database does not yet contain the signature of the new threat. To solve the problem of detecting zero-day attacks, AIDS is used.

Because it can surpass SIDS' limitations, AIDS has attracted much interest from researchers. AIDS creates a computer system behavior model based on machine learning, statistical, or knowledge-based techniques. An anomaly, which is viewed as an intrusion, is any significant divergence from the model. This set of methods is predicated on malicious conduct deviating from user behavior. AIDS can identify unknown or zero-day attacks because It is independent of signature databases. Such as SIDS to detect abnormal behavior.

IDS divided the data into two groups depending on the input source: Host IDS(HIDS) and Networks IDS (NIDS). HIDS analyzes data from the host system, including sources such as operating system logs, firewall logs, and database logs. IDS can recognize insider attacks that don't use network traffic, where the NIDS analyzes information from sources like packet capture to keep track of network activity. It can be utilized to monitor several computers linked to a network and detect early signs of external malicious activity before it spreads to other systems. To identify abnormal activity, ML algorithms, including XGBoost, TD, RF, Gaussian NB, LR, and SVM, have been used in the AIDS domain. The following section explains the ML algorithms used in this paper.

*C. ML Classifiers*

ML is a technique for instructing computers to learn from data without explicit programming. It is a subfield of artificial intelligence that enables systems to improve automatically over time. Supervised, unsupervised, and reinforcement learning are the three main types of machine learning [23, 24]. Supervised learning solves issues when the only available data consists of labeled instances. Unsupervised learning is used to find the pattern of unlabeled data. With reinforcement learning, a computer agent learns to perform a task by repeatedly attempting it and modifying its behavior in response to the feedback it receives. The model that is trained to categorize input data into classes or categories is known as a classifier. There are numerous classifier varieties, each with unique advantages and disadvantages. The data's properties and the problem's nature determine which classifier should be used. For example, while some classifiers perform better when given high-dimensional data, others perform better with fewer features. Furthermore, although some classifiers are more susceptible to noise or outliers in the data, others are more

resistant. This study selects the following Supervised learning classifiers to detect the DoS attack in IoT Environments because they have been extensively utilized in previous research on the detection of DoS attacks.

- XGBoost

XGBoost stands for Extreme Gradient Boosting, extending a version of gradient boosting [25]. It is solid and compelling, handling big datasets and producing reliable predictions. The XGBoost operates by creating a model made up of several decision trees. It begins by using the input data to train a straightforward decision tree and iteratively adds new decision trees to the model while fixing the flaws in the earlier trees. A more potent and precise model is produced due to this procedure, which is referred to as boosting. Additionally, the XGBoost provides many features, including support for missing values, management of categorical variables, and handling of imbalanced data. Additionally, it has built-in regularization to avoid overfitting and supports parallel processing to quicken the training process.

- DT

DT is a supervised ML classifier that can be used to solve classification issues [26]. It builds a tree-like model of choices and their potential effects, with each internal node standing in for a feature or attribute and each leaf node for a class label. The method recursively separates the data according to the feature values starting at the root node until it reaches a leaf node. This leaf node's class label is then used as the anticipated class for the incoming data. Decision trees are straightforward to use, easy to grasp, and capable of handling category and numerical data. But if the tree is too deep or complicated, it could be prone to overfitting.

- RF

An RF is a kind of ensemble learning classifier for classification and regression that builds several DTs during training. It produces the class that represents the mean of the classes (classification) or mean prediction (regression) of the individual trees [26]. The model performs better overall and has less overfitting when numerous decision trees are created, as opposed to depending just on one. The name's "random" component alludes to the randomly selected subsets of data that were used to train each DT.

- LR

LR attempts to calculate the chance of a particular outcome given a specific input variable. This outcome is generally binary, meaning it is composed of two possible values, such as true or false, yes or no, and so on. Multinomial logistic regression can be used to address regression with more than two possible outcomes. Logistic regression is instrumental in discovering which group a novel sample is most like. Furthermore, it is beneficial in cyber security since most security issues are categorization problems, such as recognizing attacks.

- SVM

SVM, a supervised learning technique, can be applied to classification and regression problems [27]. The fundamental

goal of SVM is to identify the optimal boundary (or hyperplane) for classifying the data into multiple groups. A boundary's margin, or the distance between it and the nearest data points from each class, should be maximized to be considered the optimal boundary. Support vectors refer to these nearby data points. By determining which side of the border additional data points fall on once the boundary has been identified, it is simple to classify them. By translating the data into a higher dimension where it may be linearly separable, SVM can also be utilized for data that cannot be separated linearly.

*D. DoS Attack*

A DoS is a type of cyberattack [28] in which the attacker tries to prevent the targeted users from using a computer resource by flooding it with a lot of traffic or requests. For example, a program, network, or website may be the target of a DoS attack to deny access to legitimate users. DoS attacks come in a variety of forms, including:

- Flooding attacks, such as network flood attacks that saturate networks with large packets, overload a targeted resource with traffic.

- Amplification attacks, increase the amount of traffic directed at a targeted resource by making use of a flaw in it, such as a Domain Name System (DNS) amplification attack that makes use of a DNS server to increase traffic to a targeted resource.

- Application-layer attacks, such as a Hypertext Transfer Protocol (HTTP) request flood attack that floods a targeted website with HTTP requests, target certain flaws in an application.

Using firewalls, IDS, machine learning, and other security measures can assist in detecting and preventing DoS attacks. The following section will review research that researchers did to develop IDS that could detect DoS attacks in the IoT using machine learning.

## III. RELATED WORKS

This section investigates previous studies that are related to the works of this paper. Much research has been done using ML to detect DoS attacks in IoT environments. Here are a few of them that will be presented.

In a study by Shreekhand and Deepak [29], in this study, they proposed and implemented ML and neural network models based on Multilayer Perceptron (MLP) and RF for the detection of DoS attacks. The proposed model successfully identified application-layer DoS attacks. The findings indicate that, compared to the MLP algorithm, which offers an accuracy of 98.87%, the RF algorithm provides a better accuracy of 99.95%. The proposed model was examined with the CICIDS2017 dataset.

A study by Yasin. et al. [30], in this study, various ML classifiers have been examined for identifying various DDoS attacks. k-Nearest Neighbors (KNN), LR, MLP, naïve Bayes, SVM, RF, deep autoencoder, CatBoost, Stacking, and XGBoost classifiers are among the ML classifiers mentioned above. The most accurate classifiers are stacking random forest and catBoost. In addition, the proposed model has been examined with the Labris and Digiturk datasets.

A study by Naeem et al. [31], in this study, a model for detecting DoS attacks during Message Queuing Telemetry Transport (MQTT) attacks in the IoT is proposed. Averaged one-dependence estimators (AODE), C4.5, and MLP ML classifiers were used to validate the proposed model. The results show that the AODE classifier achieved the best classification accuracy in identifying the DOS attack.

Another study by Jiyeon Kim et al. [32], in this study, propose a Convolutional Neural Network (CNN) based algorithm for identifying DoS attacks by taking into account the size of the kernel and the number of convolutional layers and then comparing the proposed model with Recurrent Neural Networks (RNN). The proposed model has been examined with KDD and CSE-CIC-IDS 2018 datasets. The obtained results show that the CNN outperformed the RNN.

A study by Rios. In this study, et al. [33] proposes a model for detecting DOS attacks in communication networks based on a Broad Learning System (BLS) that produces high results with less time training. The proposed model has been examined with CICIDS2017 and CSE-CIC-IDS 2018 datasets. The results show that Non-incremental BLS frequently produced the highest accuracy and F-Score, although BLS with incremental learning typically required less training time.

In another study by Verma and Ranga [34], in this study, an extensive investigation into anomaly-based IDS for protecting IoT from DoS attacks is conducted for seven ML classifiers. These classifiers include the RF, Adaboost, Gradient Boosted Machine, XGBoost, Extremely Randomized Trees, classification and regression trees, and MLP. The investigation model has been examined with frequently used datasets, including CIDDS-001, UNSWNB15, and NSL-KDD. The statistical analysis of performance metrics is conducted using Friedman and Nemenyi post-host tests to identify significant differences between classifiers. The results show that classification, regression trees, and the XGboost classifier have the best response time.

In a study by Muhammad Zeeshan et al. [35], in this study, protocol-based deep intrusion detection is proposed for IoT networks. The proposed protocol aims to find similar features of UNSW-NB15 and the Bot-IoT by comparing them. The outcome of the comparison was extracting 26 features and combining normal packets from UNSW-NB15 and DoS/DDoS from Bot-IoT. Furthermore, the proposed protocol was trained using the Long short-term memory networks (LSTM) deep learning technique, and The results show that classification accuracy was 96.3%.

In a study by Alaeddine Mihoub et al. [36], in this study, a model is proposed for the IoT to identify and prevent DoS/DDoS attacks. The identification part of the proposed model is based on the multi-class classifier that adopts the "Looking-Back" concept. Used ML classifiers in the model are DT, RF, KNN, MLP, LSTM, and RNN. The model was tested and evaluated using the Bot-IoT dataset. The obtained results show that the Looking-Back-enabled RF classifier achieves the best accuracy.

In another study by Alimi et al. [37], in this study, a model is proposed for detecting DoS attacks in IoT based on a redefined LSTM deep learning approach. The model was tested and evaluated using NSL-KDD and CICIDS-2017 datasets. The conducted results show that the proposed model has a detection accuracy of 99.22% for DoS attacks on the CICIDS-2017. In comparison, the NSL-KDD dataset attained 98.60%.

In the study by Kimmi Kumari and M. Mrunalini [38], in this study, mathematical and ML model models have been proposed for DoS attack detection using Logistic Regression and Naive Bayes. The proposed models have been tested and evaluated using the CAIDA dataset 2007. The obtained results show the mathematical model is 99.75% accurate, while the ML model is 100% accurate.

## IV. PROPOSED MODEL

Detailed step-by-step instructions of the proposed model and an explanation are provided in this section. Fig. 3 shows the proposed model flowchart.



Fig. 3. Proposed model flowchart.

### A. UNSW-NB15 Dataset

The University of New South Wales in Sydney, Australia, developed the network intrusion detection dataset known as UNSW-NB15 [39]. It is made to aid intrusion detection research and serve as a realistic testbed for evaluating IDS. It is frequently used for research and development of IDS because it comprises many actual network attacks as well as normal traffic. As a result, the dataset is commonly used in academic research projects and publications by the cybersecurity research community. The UNSW-NB15 dataset includes nine different kinds of network attacks, including worm attacks, backdoor attacks, DoS attacks, exploits attacks, fuzzers attacks, generic attacks, reconnaissance attacks, shellcode attacks, and generic attacks. Fig. 4 shows the attack distribution of the UNSW-NB15 dataset.



Fig. 4. Attack in the UNSW-NB15.

This paper focuses on the DoS attack; therefore, the DoS attack and normal traffic have been extracted from the UNSW-NB15 dataset.

### B. Dataset Preprocessing

The following preprocessing procedures must be done on the UNSW-NB15 dataset before it can be used with the proposed models.

- Label Encoding: In this stage, the category variables are transformed into numerical form so that algorithms can interpret them. This encoding process is accomplished by giving each category in the dataset an individual number. This is helpful since many ML algorithms perform better with numerical data than category data. Reduced data dimensions and enhanced model performance are two benefits of label encoding. The Label Encoding is done using The Scikit-Learn Library's preprocessing module in Python.

- Normalization: Normalization is a technique used in dataset processing by transforming and scaling data to fit inside a predetermined range. The purpose of doing this is often to lessen the influence of outliers and guarantee that the data falls within a close range. The Min-Max normalizing method is a popular approach where the data is often between 0 and 1. The following formula can be used to normalize a value x using the Min-Max normalization:

$$\text{Min/Max normalization} = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (1)$$

- Remove the missing: In a dataset, removing the missing value is a part of the cleaning process that eliminates errors, inconsistencies, and unwanted information. This process ensures that the data is accurate, consistent, and prepared for analysis easier. There are several ways to remove missing values (NaN values) from a dataset using Python. The dropna() method from the Pandas library and the fillna() method to replace a given value for missing values were used.

### C. Dataset Splits

Dataset splitting is breaking up a large dataset into more manageable chunks for uses like model testing and training. A dataset is frequently divided into Training sets: this is the primary dataset used to train a model. It is utilized to discover the underlying patterns in the information and contains a sizable amount of data. Testing set: The test set is used to evaluate how well the trained model has worked in practice. It contains data, the model hasn't seen before and estimates how well it performs on actual data.

The split ratios for training and test sets vary depending on the size and complexity of the dataset, but a typical split ratio is 70–30. (Training - Test). Therefore, ensuring the split is accurate and the various subgroups do not overlap is crucial. Python's most used data split method is train_test_split from the scikit-learn library.

### D. ML Classifiers

An ML classifier is used to categorize incoming data as either abnormal or normal. The XGBoost, DT, RF, NB, LR, and SVM classifiers are discussed in detail in Section II(C). Because these are the most well-known classifiers used in the literature for IDS, these classifiers were chosen.

### E. Confusion Matrix and Model Evaluations

In machine learning, a confusion matrix is a table that assesses how well a classification model performs. It lists the model's accurate and inaccurate predictions based on a test data sample. True Positives (TP), False Positives (FP), True Negatives (TN), and False Negatives (FN) are the confusion matrix's four main components. These components are used to produce some metrics, including F1-score, recall/sensitivity, accuracy, and precision, which provide a comprehensive picture of the model's performance.

$$\text{Accuracy}=(TP+TN)/(TP + TN +FP +FN) \quad (2)$$

$$\text{Precision}= TP /(TP + FP) \quad (3)$$

$$\text{Recall}= TP /(TP + FN) \quad (4)$$

$$\text{F-Measure} =(2 \times \text{Precision} \times \text{Recall} )/(\text{Precision}+\text{Recall}) \quad (5)$$

### V. EXPERIMENTAL DESIGN AND RESULTS

This section presents the Experiment Design and Results

### A. Experimental Design

The model was evaluated using Windows 7 and an i7 processor running at 3.40 GHz with 6.0 GB of RAM. The experiments were conducted using the open-source Anaconda (spider) for the UNSW-NB15 datasets under both normal and DoS attacks. Scikit-learn tools in Python were used to implement the model. This model uses the classifiers DT, RF, NB, LR, and SVM. The obtained confusion matrix of prediction for several classifiers is shown in Fig. 5.



Fig. 5.  Confusion matrix of each classifier.

Table I presents the outcomes of the model's assessment metrics by employing the confusion matrix, which was attained, as depicted in Fig. 7.

TABLE I.     PERFORMANCE METRICS RESULTS

| | Accuracy | Precision | Sensitivity | F-measure | TP | FN | FP | TN |
|---|---|---|---|---|---|---|---|---|
| **XG Boost** | 98.9 | 98.4 | 99.2 | 98.8 | 99.2 | 0.7 | 1.4 | 98.5 |
| **SVM** | 98.6 | 98.2 | 98.7 | 98.5 | 98.7 | 1.2 | 1.5 | 98.4 |
| **NB** | 83.7 | 77.5 | 91.8 | 84.0 | 91.8 | 8.1 | 23.3 | 76.6 |
| **LR** | 97.7 | 97.1 | 98.0 | 97.6 | 98.0 | 1.9 | 2.5 | 97.4 |
| **RF** | 99.4 | 99.2 | 99.6 | 99.4 | 99.6 | 0.3 | 0.7 | 99.3 |
| **DT** | 99.1 | 99.1 | 99.0 | 99.1 | 99.0 | 0.9 | 0.7 | 99.2 |

## B. Finding

The obtained results of this study are analyzed in this section. As seen in Fig. 6, the accuracy of the various ML classifiers is demonstrated. Accuracy is the percentage of accurate predictions made by a classifier in comparison to the actual value of the label. The accuracy rate for XGboost was 98.92%, SVM 98.62%, NB 83.75%, LR 97.74%, RF 99.48%, and DT 99.16%. According to the data collected, the RF classifier outperformed the other classifiers in terms of accuracy due to the following, the RF comprises Multiple decision trees, therefore, it has less classification error. A measure of precision is a percentage that shows what proportion of the objects the classifier recognized are accurate forecasts. For example, the precision for XGboost, SVM, NB, LR, RF, and DT in Fig. 7 is 98.40%, 98.29%, 77.50%, 97.14%, 99.21%, and 99.12%, respectively. The outcomes show that the RF classifier was more precise than the other classifiers.

Sensitivity, also known as recall or true positive rate, is a commonly used metric in machine learning to evaluate binary classification models. It indicates the number of true positive cases the classifier correctly categorized as positive. Sensitivity values for XGBoost, SVM, NB, LR, RF, and DT as in Fig. 8, 99.29%, 98.76%, 91.87%, 98.06%, 99.69%, and 99.08%, respectively. The obtained Sensitivity values demonstrate the superiority of RF classifiers over other classifiers. A classifier's performance can be assessed using the F-measure since it considers both the precision and sensitivity values. This metric is beneficial when there is an uneven distribution of positive and negative classifications. F-measure values for XGBoost, SVM, NB, LR, RF, and DT as in Fig 9, 98.85%, 98.53%, 84.08%, 97.60%, 99.45%, and 99.10%, respectively. The results show that RF classifiers outperform other classifiers regarding F-measure values.

From the analysis mentioned above, the following conclusions can be drawn:

*1) RF* is the best classifier to detect the DoS attack in an IoT environment compared to XGBoost, DT, NB, LR, and SVM.

*2) Gaussian* NB is the worst classifier to detect the DoS attack in an IoT environment.



Fig. 7.   Precision of ML classifiers.



Fig. 8.   Sensitivity of ML classifiers.



Fig. 9.   F-Measure of ML classifiers.



Fig. 6.   Accuracy of ML classifiers.

## VI. CONCLUSION AND FUTURE WORKS

Security measures must be included in IoT environments to prevent and combat DoS attacks. Therefore, this paper introduced an IDS model using XGBoost, DT, RF, NB, LR, and SVM to detect DoS attacks in the IoT environment on the UNSW-NB15 datasets. The outcomes of the model were examined using the confusion matrix. Accuracy, precision, sensitivity, and F-Measure were used to evaluate the model's performance. The obtained experimental data proves that the RF is the most efficient classifier among other examined classifiers to detect DoS attacks in IoT environments. Its accuracy was 99.48%, precision 99.21%, sensitivity 99.69%, and F-Measure 99.45%. This study was limited to evaluating some of the ML classifiers with specific attacks. In the future direction of the research, the modern dataset for IDS and other types of DoS attacks, such as (DDoS) attacks, deep learning, and reinforcement learning approaches will be considered to examine the model performance.

## REFERENCES

[1] L. Farhan, S. T. Shukur, A. E. Alissa, M. Alrweg, U. Raza, and R. Kharel, "A survey on the challenges and opportunities of the Internet of Things (IoT)," in 2017 Eleventh International Conference on Sensing Technology (ICST), 2017, pp. 1-5: IEEE.

[2] M. A. Ferrag, O. Friha, D. Hamouda, L. Maglaras, and H. Janicke, "Edge-IIoTset: A new comprehensive realistic cyber security dataset of IoT and IIoT applications for centralized and federated learning," IEEE Access, vol. 10, pp. 40281-40306, 2022.

[3] O. Almomani, M. A. Almaiah, M. MADI, A. Alsaaidah, M. A. Almomani, and S. Smadi, "Reconnaissance attack detection via boosting machine learning classifiers," in AIP Conference Proceedings, 2023, vol. 2979, no. 1: AIP Publishing.

[4] A. Almomani et al., "Ensemble-Based Approach for Efficient Intrusion Detection in Network Traffic," vol. 37, no. 2, 2023.

[5] G. Carl, G. Kesidis, R. R. Brooks, and S. Rai, "Denial-of-service attack-detection techniques," IEEE Internet computing, vol. 10, no. 1, pp. 82-89, 2006.

[6] R. Vishwakarma and A. K. J. T. s. Jain, "A survey of DDoS attacking techniques and defence mechanisms in the IoT network," vol. 73, no. 1, pp. 3-25, 2020.

[7] P. Kumari, A. K. J. C. Jain, and Security, "A Comprehensive Study of DDoS Attacks over IoT Network and Their Countermeasures," p. 103096, 2023.

[8] X. Zhu and H. Deng, "A security situation awareness approach for iot software chain based on markov game model," 2022.

[9] O. Almomani, M. A. Almaiah, A. Alsaaidah, S. Smadi, A. H. Mohammad, and A. Althunibat, "Machine learning classifiers for network intrusion detection system: comparative study," in 2021 International Conference on Information Technology (ICIT), 2021, pp. 440-445: IEEE.

[10] A. H. Mohammad, T. Alwada'n, O. Almomani, S. Smadi, N. ElOmari, and Continua, "Bio-inspired Hybrid Feature Selection Model for Intrusion Detection," Computers, Materials, vol. 73, no. 1, pp. 133-150, 2022.

[11] A. Alzaqebah, I. Aljarah, O. Al-Kadi, and R. Damaševičius, "A modified grey wolf optimization algorithm for an intrusion detection system," Mathematics, vol. 10, no. 6, p. 999, 2022.

[12] O. Almomani, "A hybrid model using bio-inspired metaheuristic algorithms for network intrusion detection system," Comput. Mater. Contin, vol. 68, no. 1, pp. 409-429, 2021.

[13] O. Almomani, "A feature selection model for network intrusion detection system based on PSO, GWO, FFA and GA algorithms," Symmetry, vol. 12, no. 6, p. 1046, 2020.

[14] P. Gokhale, O. Bhat, S. Bhat, and Technology, "Introduction to IOT," International Advanced Research Journal in Science, Engineering, vol. 5, no. 1, pp. 41-44, 2018.

[15] P. Matta, B. Pant, and Technology, "Internet of things: Genesis, challenges and applications," Journal of Engineering Science, vol. 14, no. 3, pp. 1717-1750, 2019.

[16] A. Al Zaqebah, O. Almomani, M. Almomani, A. Alsaaidah, A. A. Abu-Shareha, and A. Althunibat, "Improving Routing Decision Algorithm for RPL Networks," in 2023 International Conference on Information Technology (ICIT), 2023, pp. 544-549: IEEE.

[17] M. Almaiah, F. Hajjej, A. Ali, M. Pasha, and O. Almomani, "An AI-Enabled Hybrid Lightweight Authentication Model for Digital Healthcare Using Industrial Internet of Things Cyber-Physical Systems," Sensors, vol. 22, p. 1448, 2022.

[18] R. Masadeh, B. AlSaaidah, E. Masadeh, M. d. R. Al-Hadidi, and O. Almomani, "Elastic Hop Count Trickle Timer Algorithm in Internet of Things," Sustainability, vol. 14, no. 19, p. 12417, 2022.

[19] V. Hassija, V. Chamola, V. Saxena, D. Jain, P. Goyal, and B. Sikdar, "A survey on IoT security: application areas, security threats, and solution architectures," IEEE Access, vol. 7, pp. 82721-82743, 2019.

[20] R. Masadeh, O. Almomani, E. Masadeh, and R. e. Masa'deh, "Secure CoAP Application Layer Protocol for the Internet of Things Using Hermitian Curves," in The Effect of Information Technology on Business and Marketing Intelligence Systems: Springer, 2023, pp. 1869-1884.

[21] H. F. Atlam, G. B. Wills, and s. cities, "IoT security, privacy, safety and ethics," Digital twin technologies, pp. 123-149, 2020.

[22] D. Mudzingwa and R. Agrawal, "A study of methodologies used in intrusion detection and prevention systems (IDPS)," in 2012 Proceedings of IEEE Southeastcon, 2012, pp. 1-6: IEEE.

[23] A. Haldorai, A. Ramu, and M. Suriya, "Organization internet of things (IoTs): supervised, unsupervised, and reinforcement learning," in Business Intelligence for Enterprise Internet of Things: Springer, 2020, pp. 27-53.

[24] A. Sholiyi, J. A. Alzubi, O. A. Alzubi, O. Almomani, and T. O'Farrell, "Near capacity irregular turbo code," arXiv preprint arXiv:.01358, 2016.

[25] S. Smadia, O. Almomanib, A. Mohammadc, M. Alauthmand, and A. Saaidahe, "VPN Encrypted Traffic classification using XGBoost," International Journal of Advanced Trends in Computer Science and Engineering, vol. 9, no. 7, 2021.

[26] M. Madi, F. Jarghon, Y. Fazea, O. Almomani, A. Saaidah, and C. Sciences, "Comparative analysis of classification techniques for network fault management," Turkish Journal of Electrical Engineering, vol. 28, no. 3, pp. 1442-1457, 2020.

[27] M. A. Almaiah et al., "Performance Investigation of Principal Component Analysis for Intrusion Detection System Using Different Support Vector Machine Kernels," Electronics, vol. 11, no. 21, p. 3571, 2022.

[28] S. SMADI, M. ALAUTHMAN, O. ALMOMANI, A. SAAIDAH, and F. ALZOBI, "Application layer denial of services attack detection based on stacknet," International Journal of Advanced Trends in Computer Science and Engineering, vol. 3929, no. 3936, pp. 2278-3091, 2020.

[29] S. Wankhede and D. Kshirsagar, "DoS attack detection using machine learning and neural network," in Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), 2018, pp. 1-5: IEEE.

[30] Y. Gormez, Z. Aydin, R. Karademir, and V. C. Gungor, "A deep learning approach with Bayesian optimization and ensemble classifiers for detecting denial of service attacks," International Journal of Communication Systems, vol. 33, no. 11, p. e4401, 2020.

[31] N. F. Syed, Z. Baig, A. Ibrahim, and C. Valli, "Denial of service attack detection through machine learning for the IoT," Journal of Information Telecommunication, vol. 4, no. 4, pp. 482-503, 2020.

[32] J. Kim, J. Kim, H. Kim, M. Shim, and E. Choi, "CNN-based network intrusion detection against denial-of-service attacks," Electronics, vol. 9, no. 6, p. 916, 2020.

[33] A. L. G. Rios, Z. Li, K. Bekshentayeva, and L. Trajković, "Detection of denial of service attacks in communication networks," in IEEE

international symposium on circuits and systems (ISCAS), 2020, pp. 1-5: IEEE.

[34] A. Verma and V. Ranga, "Machine learning based intrusion detection systems for IoT applications," Wireless Personal Communications, vol. 111, pp. 2287-2310, 2020.

[35] M. Zeeshan et al., "Protocol-based deep intrusion detection for dos and ddos attacks using unsw-nb15 and bot-iot data-sets," IEEE Access, vol. 10, pp. 2269-2283, 2021.

[36] A. Mihoub, O. B. Fredj, O. Cheikhrouhou, A. Derhab, and M. Krichen, "Denial of service attack detection and mitigation for internet of things using looking-back-enabled machine learning techniques," Computers Electrical Engineering, vol. 98, p. 107716, 2022.

[37] K. O. Adefemi Alimi, K. Ouahada, A. M. Abu-Mahfouz, S. Rimer, and O. A. Alimi, "Refined LSTM Based Intrusion Detection for Denial-of-Service Attack in Internet of Things," Journal of Sensor Actuator Networks, vol. 11, no. 3, p. 32, 2022.

[38] K. Kumari and M. Mrunalini, "Detecting Denial of Service attacks using machine learning algorithms," Journal of Big Data, vol. 9, no. 1, pp. 1-17, 2022.

[39] N. Moustafa and J. Slay, "UNSW-NB15: a comprehensive data set for network intrusion detection systems (UNSW-NB15 network data set)," in 2015 military communications and information systems conference (MilCIS), 2015, pp. 1-6: IEEE.

# Improving the Trajectory Clustering using Meta-Heuristic Algorithms

Haiyang Li*, Xinliu Diao

Department of Basic Courses, Henan Polytechnic Institute, Nanyang Henan 473000, China

*Abstract*—**The rapid growth of GPS trajectories obscures valuable information regarding urban road infrastructure, urban traffic patterns, and population mobility. An innovative method termed trajectory regression clustering is introduced to improve the extraction of hidden data and generate more precise clustering results. This approach belongs to the unsupervised trajectory clustering category and has the objective of minimizing the loss of local information inside the trajectory. It also seeks to prevent the algorithm from getting stuck in a suboptimal solution. The methodology we employ consists of three primary stages. To begin with, we present the notion of trajectory clustering and devise a distinctive approach known as angle-based partitioning to segment line segments. The evaluation results indicate a significant improvement in the clustering accuracy of the proposed method compared to existing methodologies, especially for a high number of clusters. The HCMGA and HCMMOPSO algorithms have improved clustering accuracy for MBP values by 0.61% and 0.64%, respectively, as compared to previous approaches. Moreover, based on the implementation findings, the ant colony approach demonstrates superior accuracy compared to alternative methods, while the particle swarm method exhibits faster convergence.**

*Keywords*—*Ant colony method; particle swarm algorithm; HCM clustering; and trajectory lines*

## I. INTRODUCTION

An analysis of the travel characteristics of moving objects, such as automobiles and individuals, can provide insights into travel patterns. This analysis reveals information about people's frequent travel habits, patterns of traffic congestion, and social activity patterns [1]. Travel patterns have been utilized in various domains, such as furnishing decision-making data for urban planning and emergency situations [2], analyzing and optimizing routes to offer personalized travel suggestions for residents, dispatching vehicles [3], and optimizing and selecting stations [4]. These applications can provide valuable insights for urban construction and growth. The authors in Reference [5] proposed a method that takes into account the road network while clustering road segment spatial trajectories. This method was designed to replace density-based clustering and Euclidean-based distance calculations, with the goal of achieving faster and more efficient clustering. The algorithm described in reference [6] is a scalable and efficient density clustering method that utilizes big data computing. In addition, the authors of Reference [7] introduced an enhanced density-based technique specifically designed for clustering stops in trajectories. The study described in Reference [8] introduced a clustering approach based on anisotropic density (angle-based

standard deviation). This algorithm was utilized to identify spatial point patterns along with any accompanying noise.

Typically, clustering techniques can be classified into many categories such as density-based, partitioning-based, grid-based, hierarchical-based, and graph-based. These algorithms are extensively used in spatial data processing for a variety of purposes. In addition, each of these categories encompasses other renowned clustering algorithms, such as partitioning-based K-means, K-median, and fuzzy C-means (HCM), each with its own distinct advantages and disadvantages. Specifically, density-based clustering algorithms are commonly employed to extract concealed information from a given dataset and process any GPS datasets, since they are especially well-suited for identifying clusters with irregular forms and identifying clusters that do not overlap [9, 10, 11]. Nevertheless, managing the intersecting clusters (such as trajectory crossover) becomes challenging when dealing with fuzzy clusters and the absence of localized trajectory information. Furthermore, it is responsive to both the specified neighborhood and the density, as determined by the value of MinPts. This research specifically examines partitioning-based techniques, such as HCM. Nevertheless, they exhibit certain limitations, such as being susceptible to the choice of initial cluster centers, delayed convergence, and a propensity to get trapped in local optima. In this paper, a new trajectory regression clustering technique is introduced. The technique is based on partition clustering and combines the AngPart method, which generates line segments based on angles, with the HCML algorithm, a Lagrange-based fuzzy C-means clustering algorithm, and the LSR model, which is used to create an unsupervised trajectory clustering method. This method is an alternative to using a map-based knowledge base. Specifically, HCML is an innovative approach for clustering regression data without the need for supervision. Initially, a line segment partitioning technique is developed to generate line segments using three GPS data points. This method efficiently preserves the local information of trajectories. The novel clustering algorithm presented in this study combines a unique fuzzy C-means (NHCM) with the Lagrange operator [12] and Hausdorff-based K-means++ [13]. The NHCM is employed to cluster line segments, while K-means++ is utilized to generate the initial cluster centers of the line segments. This combination aims to capture the global optimum and prevent the algorithm from getting trapped in local optima. The original fuzzy C-means (HCM) technique is a clustering approach based on partitioning [14]. In this algorithm, the computation of distances between line segments requires the use of Hausdorff distances instead of Euclidean distance [15]. Once the concealed GPS data is extracted and

acquired, the LSR (Least Squares Regression) method is utilized to perform trajectory regression. The objective is to regress and generate trajectories based on the clustering results, without relying on a map-based knowledge base. These trajectories can then be used to analyze and describe various urban aspects such as people, vehicles, roads, traffic flow, and serve as a reference for road planning.

*Background:* HCML can enhance line segment partitioning and maintain the local information of trajectory by employing the angle-based technique prior to the clustering process. For instance, when two GPS data points are produced as line segments, it becomes challenging to articulate the local information between GPS data points and to comprehend the connection between consecutive GPS points, such as the alteration in steering and intersecting angle.

Furthermore, the method being discussed, HCML, is a form of unsupervised learning. Hence, in cases when a map-based knowledge base is not required, the least squares regression model (LSR) is employed to generate the trajectories of the clustering outcomes.

*Problem:* In order to assess the performance and efficacy of HCML, an actual GPS dataset from Beijing, China is employed as an experimental test. The experiments involve comparing HCML with K-median, K-means, and HCM clustering methods using the MBP (Pakhira-Bandyopadhyay-Maulik)-index [26] cluster evaluation criteria. This is discussed in Section 2. The MBP-index is an effective unsupervised evaluation tool [27,28]. However, it is important to be aware that distances in the MBP-index necessitate the utilization of the Hausdorff method for calculating the distance between the center of line segments and other line segments. Furthermore, LSR is employed to accomplish the regression of the clustering outcomes. The experimental findings demonstrate that HCML outperforms K-means, K-median, and HCM algorithms in terms of trajectory regression quality (refer to Section 5).

Proposed Solution: Thus, the primary content of the study is succinctly outlined as follows:

*1)* The angle-based partitioning method (AngPart) is offered as a strategy for generating unique line segments.

*2)* This study proposes a novel clustering technique called fuzzy C-means (NHCM), which integrates the Lagrange operator with AngPart and K-means++.

*3)* This paper introduces a trajectory regression technique that utilizes least squares regression (LSR) to analyze population movement patterns along trajectories. The technique provides insights into the state of population migration and can serve as a valuable reference for urban road planning.

*4)* HCML has been demonstrated to be effective when used to actual taxi GPS data in Beijing, China.

The subsequent sections of the paper are arranged in the following manner. Section II provides a description of taxi GPS data in Beijing, China. Section III presents the angle-based normalization method employed for the taxi GPS data. Section IV introduces a trajectory regression technique that combines HCML with LSR. Section V outlines the tests and offers the data for evaluating the effectiveness of the recommended procedures. Section VI serves as the final part of the report, providing a conclusion and proposing potential future research.

## II.    RELATED WORKS

The integration of information technology into transportation systems is currently a prominent trend. This is because it effectively addresses key issues faced by traffic operators, such as traffic congestion and accidents. Therefore, observing the traffic situation is essential for traffic operators, particularly at intersections [16]. The traffic data obtained from the monitoring system is frequently extensive, necessitating diligent attempts to identify noteworthy trends within it. These patterns provide valuable insights into vehicle movements and facilitate the detection of any deviant conduct that may result in traffic disputes. Nevertheless, it will be an arduous task for traffic operators to manually monitor the movement of vehicles at a crossroads, especially when there are thousands of vehicles passing through. Therefore, the process of grouping vehicle trajectory data to identify comparable patterns is carried out using the k-means and fuzzy c-means (HCM) clustering algorithms. Since various clustering techniques necessitate the input parameter of the number of clusters, this research focuses on studying the appropriate number of clusters for the clustering process [17].

Analyzing urban travel patterns helps assess the regularity of inhabitants' mobility, offering information for urban traffic planning and emergency decision-making. Clustering techniques have been extensively utilized to uncover concealed insights from extensive trajectory data concerning trip patterns. Implementing soft constraints in the clustering process and statistically evaluating their performance remains a challenging task. This paper introduces a refined trajectory clustering approach, known as TC-FDBSCAN, which utilizes fuzzy density-based spatial clustering of applications with noise to perform classification on trajectory data. Initially, we establish the trajectory distance by taking into account various features and determining the corresponding weight factors to quantify the similarity between trajectories [18]. In the HCM clustering method, membership degrees and membership functions are created to extend the standard DBSCAN method.

A modified version of the fuzzy c-means algorithm is proposed to solve the inverse kinematics and trajectory planning problem for a redundant manipulator, while taking into account performance criteria. A novel HCM clustering approach is introduced, which utilizes a newly developed generalized validity index. This index is built on weighted within-scatter metrics and between-cluster scatter metrics specifically designed for the manipulator. The issue of redundant manipulator, which refers to a nonlinear system with several inputs and outputs, has not been previously addressed using the clustering method. The trajectory planning algorithm for the manipulator is simulated using Matlab. The entire process, starting from gathering data to verifying the model, is demonstrated using a robot manipulator with four degrees of freedom. The simulated results are being compared to the numerical approaches used for trajectory planning. The findings are visually depicted. The proposed method offers the

benefits of being straightforward, adaptable, and exhibiting excellent tracking capabilities [19].

The mining of trajectory databases (TD) has garnered significant attention as a result of the widespread use of tracking devices. However, the mining procedure has not yet included the presence of uncertainty in TD, such as GPS inaccuracies. This work examines the impact of uncertainty in TD clustering and presents a three-step methodology to address it. Initially, we give a conceptual framework for representing trajectories using intuitionistic point vectors, which captures the inherent uncertainty. Additionally, we introduce a robust distance metric to handle this uncertainty. Furthermore, we provide CenTra, an innovative approach that addresses the challenge of identifying the centroid trajectory of a set of moves. Furthermore, we provide a modified version of the fuzzy C-means (HCM) clustering algorithm that incorporates CenTra during its updating step. The empirical assessment on real-world TD substantiates the efficiency and efficacy of our methodology [20].

This research examines the application of HCM clustering to identify probable spatial patterns by incorporating rough set and fuzzy set theory. Initially, we suggest a rapid technique for measuring similarity by utilizing the approximate distances between trajectories. Significant reductions in processing time would be achieved, particularly for lengthy trajectory sequences. In addition, we present a summarization strategy that minimizes the number of distance calculations needed for similarity measurement. Furthermore, the membership degree is modified in order to enhance the clustering quality and performance. Furthermore, we enhance the fuzzy C-means algorithm by using a novel similarity measure and the membership degree function. The efficacy of our methods is demonstrated by experimental findings obtained from two actual datasets of trajectories. These results involve the assessment of clustering validity and the computation performance for big datasets. The computing performance of the proposed HCM clustering algorithm exhibits a clear improvement as the size of the trajectory dataset rises [21].

Clustering trajectory data is a method used to identify and display the underlying structure in the movement patterns of mobile objects. This technique has a wide range of possible applications in fields such as traffic control, urban planning, astronomy, and animal research. This study introduces a novel method for grouping trajectory data using a Particle Swarm Optimization (PSO) methodology. The strategy takes into account the Dynamic Time Warping (DTW) distance, which is a widely-used measure for comparing trajectory data [22]. The suggested technique may identify the (almost) optimal number of clusters and the (almost) ideal cluster centers during the clustering process. In order to enhance the performance of the suggested method, a Discrete Cosine Transform (DCT) representation of cluster centers is utilized. This approach helps to minimize the dimensionality of the search space and improve the method's performance in relation to a certain performance index. The suggested method can incorporate

different cluster validity indices as the objective function for optimization. The experimental findings, obtained from both synthetic and real-world datasets, demonstrate the improved performance of the proposed technique compared to fuzzy C-means, fuzzy K-medoids, and two evolutionary-based clustering techniques previously suggested in the literature [23].

Next, the NT algorithm utilizes the characteristics of noise in order to actively reduce the impact of noise [24]. In 2022, a method was provided for grouping ship trajectories at sea. This approach utilized the Douglas-Peucker compression technique and the DBSCAN algorithm. The arrangement of this data defines the dispersion of traffic volume and the customary path for ship passage. Initially, the appropriate parameters are derived for the Douglas Poker method by analyzing the alterations in the ships' trajectories. This is done with the aim of enhancing the clustering of the ships' trajectories. The DTW distance matrix is calculated using these parameters to compress the data obtained from the traffic lines. Next, the enhanced DBSCAN algorithm is utilized to accomplish density-based clustering. The DBSCAN algorithm selects its optimal settings by considering the statistical properties of the distribution of ships' routes [25].

## III. Proposed Method

The suggested approach will present the path utilizing an innovative pre-processing technique. This method utilizes the angle formed between the lines as a reference to determine the linear areas in a specific order based on the sequence of points. Subsequently, an attempt is made to employ evolutionary approaches to better the performance of particle swarms and ant colonies in order to address the limitations of the HCM clustering technique. The overall framework of the suggested technique is illustrated in Fig. 1.

### A. Using GPS Coordinates to Determine Starting and Ending Places

The trajectory in the proposed approach comprises the GPS coordinates specified by the user for both people and moving vehicles. These figures represent data about a motion that starts at one point and concludes at another. During this process, sampling is conducted at regular intervals, often less than two minutes, and the spatial data of the moving item is recorded [26]. This enables the collection of data related to a sequence of locations for the mobile entity (refer to Fig. 2). The suggested technique utilizes route information represented as $Ti = \{(p, a)\} = \{(p1, a1), (p2, a2), (pi, ai)\}$, where $pi$ is a pair consisting of latitude and longitude coordinates, and $a$ represents the angle formed by the line from the current point to the next point.

### B. Tracing the Pathways of GPS Data

The results pertaining to the clustering which involves the regionalization of the trajectories are comparable. Fig. 3 illustrates an instance of the clustering of garlic lines.

Fig. 1.   General structure of the proposed method.



Fig. 2.   A GPS moving point example with an angle-based trajectory.

## C. Similarity Criteria

Among the newly-introduced criteria for comparing the similarity of linear sections is the cosine similarity criterion. Similarity results of linear regions form the basis for both trajectory clustering and sub-trajectory computing. The proposed method expresses the trajectory of Lj as a function of angle changes. The similarity criterion is also produced by relationship (1 [27].

$$sim(S_{j'}, S_j) = \begin{cases} e^{\forall l} \times \frac{e^{\partial g} - e^{-\partial g}}{e^{\partial g} + e^{-\partial g}} & S_{j'} \neq S_j \\ 1 & otger \end{cases} \quad (1)$$

The results pertaining to the clustering which involves the regionalization of the trajectories are comparable. Fig. 3 illustrates an instance of the clustering of garlic lines.

## D. Establish Zoning Lines

The process of dividing the two-part line areas is detailed here. Every possible angle for the specified GPS locations is computed and saved in normal mode when the connecting line between two consecutive points forms a right angle. The angle between the lines and the vertical line is considered. First, we find the tenth neighboring point's shortest distance from the selected point. The sums of these points can vary. This leads to the identification of a region that matches the three nearest existing places, as illustrated in Fig. 4. If, among these three locations, the largest angle measures more than 180 degrees, the guiding angle is computed counterclockwise. In cases when the maximum angle falls short of 180 degrees, the guiding angle is calculated in the opposite direction. The values associated with the angle are changed according to relations (2) and (3) to create a standard that is both normal and harmonious [28, 29, 30].

Fig. 3. Clustering of trajectories.



Fig. 4. The process of choosing three adjacent locations.

$$Nita = 2\gamma - (Nita_2 - Nita_1) \qquad (2)$$

$$Nita = Nita_2 - Nita_1 \qquad (3)$$

usage of angle-based segmentation necessitates limiting the intersection angle of the lines, denoted as $\gamma_t = (f = 1,2,\dots,e)$ The theorem of cosines, as stated in equation (4), can be used to determine the angle of the intersecting lines when three GPS points are represented as $(S^-, S, S^+)$ and P is taken as the vertex. (a) Point P is initially selected at random. (b) The vertex is selected at this position, and two nearby points are identified by calculating the lowest distance. (c) If the values of If $\gamma <$T, where T is the threshold value, are extracted from the data set and saved in memory, then step 4 is carried out; otherwise, we return to step 1. (d) In order to build linear regions, GPS points are traveled until each point is checked. Using angle-based segmentation and cosine-based limiting, linear areas are segmented in this part. Depending on the angle between them, a collection of journey lines tied to human or point device movement based on GPS points is expressed as specific lines. A greater amount of data is available for grouping by the trend lines.

### E. HCM Algorithm-based Traffic Line Clustering

In order to tackle this problem, researchers are investigating meta-heuristic algorithms. The HCM method's implementation and assessment do not effectively optimize the output of the objective function algorithm, resulting in the clustering centers being situated in local optima. Using meta-heuristic approaches improves the HCM algorithm. The defining feature of optimization and random search algorithms that rely on collective intelligence is the collective behavior and self-organization of individuals inside the community. This work

involves the grouping of trajectory data using a combination of classical and meta-heuristic clustering algorithms.

### F. Using the Particle Swarm Technique to Optimize HCM Clustering

The clustering problem is reframed as an optimization problem in the proposed strategy. In equation (4), C is a specific clustering of the data set f of the optimization function, and the optimization problem is specified as the set of all possible optimal clusterings (D*= {D1, D2,..., Dk}). C* is the best clustering that an iterative algorithm can provide [31].

$$E(D^*) = \min E(D); \qquad (4)$$

The particle swarm optimization approach is employed to get an optimal segmentation. The HCM algorithm is a clustering technique that relies on an objective function. The objective function, denoted by equation (5), is defined for the data set Y = {Y1, Y2, ..., Yn} of dimension s.

$$I_n = \sum_{i=1}^{p} \sum_{j=1}^{r} \mu_{lb}^n \|P_i - d_r\|^2 \qquad (5)$$

When determining the distance and degree of similarity of data with the center of the cluster, Let S represent the number of linear segments extracted in the preceding phases. D represents the number of clusters. m specifies the fuzzy degree of overlap of the clusters. The value of μid, which is the degree of the membership function Dj for cluster k, is determined using equation (6) [32].

$$\mu_{id} = \frac{1}{\sum_{d'=1}^{d}\left(\frac{\|P_j - E_d\|}{\|P_j - E_{d'}\|}\right)^{\frac{2}{n-1}}} \qquad (6)$$

The suggested methodology utilizes HCM clustering, employing the particle swarm optimization technique. In particle swarm optimization, the fitness function is replaced by the objective function in HCM clustering. The sequence of actions performed by the algorithm is as follows:

- Each particle in the particle swarm optimization process is assigned a membership function ujk, which is determined based on equation (6).

- Equation (7) [33] reevaluates the centroids of the categories determined using HCM clustering.

- The fitness function of the present particle is defined by Equation (5).

- The user's text is a bullet point. The ultimate outcome is the global optimum, which represents the most optimal solution when taking into account all the particles. This clustering represents the highest level of optimization and is considered the global optimum.

$$W_i(d+1) = \frac{\sum_{i=1}^{m} y_{ij}^n(d) Y_i}{\sum_{i=1}^{m} y_{ij}^n(d)} \qquad (7)$$

*G. Optimization of HCM Clustering by Ant Ga Algorithm*

As mentioned earlier, HCM clustering aims to determine the appropriate values for cluster centers and membership function values in order to achieve optimal clustering and minimize the desired objective function. Currently, there exist two separate optimization concerns. The data is subsequently divided into d clusters. The values of the pheromone matrix p are responsible for performing this task. Following each step, the pheromone matrix p undergoes modification, and the updated values for ujk are computed using equation (6) and the cluster centers. The pheromone levels are updated using Equation (8) [31].

$$s_{id} = s_{id} \times (1 - \partial) + (y_{id}/(I - I_{min} + \epsilon)^{\partial} \qquad (8)$$

for cluster d's route i, the pheromone concentration is denoted as $s_{id}$. While j is the number of samples, d is the number of clusters. The variables $\epsilon$, which precludes division by zero, $\partial$, which controls the pace of convergence, and $\partial$, which controls the amount of pheromone evaporation, are all parameters in the system. The Algorithm 1 uses an ant colony method to show the structure of HCM clustering optimization. Included in this algorithm are the input data set Y, the fuzzy power m, the number of clusters c, and the maximum number of steps needed to run the ant colony process.

---

**Algorithm 1: Optimization of HCM clustering by GA algorithm**

dataset%

Initialise Y, d, n;

Initialise Jmin = inf, **Y** = 0

Initialise GA parameters - smax, $\partial$, _, $\beta$;

Initialise the pheromone matrix, **S**, using Eq. (5);

**for** t = 1 to smax **do**

   repeat

     for k = 1 to m do

---

With probability $s_{jd} / \sum_{i=1}^{s} s_{jd}$

   set $j^{th}$ cluster membership value = 1,

   for i ≠ j, set the membership value= 0;

  end for

**until** there is at least one point per cluster;

Using Eq. (3), compute centroids **W**;

Using Eq. (2), compute the new fuzzy

membership matrix, **Y**;

Using Eq. (3), compute the new centroids **W**;

Using Eq. (1), calculate objective function I;

if I < $I_{min}$ then

   $I_{min}$ = I;

end if

Using Eq. (4), update pheromone matrix **s**;

end for

---

- Data Collection:

A large dataset consisting of GPS trajectories of cabs in operation in a big Chinese city is used in the study. All sorts of spatial and temporal dimensions are present in the dataset, which records cab travels over a long period of time. For accurate data representation, the trajectories are sampled frequently. Each trajectory point is annotated with details like timestamp, latitude, and longitude. These details are used for further research.

- Preprocessing:

The trajectory data is thoroughly preprocessed to remove noise, outliers, and missing values before grouping. In order to make the following clustering more resilient, we exclude outliers found via spatial and temporal analysis. To preserve the flow of time, values that are missing are filled in by extrapolating from nearby trajectory points. In order to simplify the dataset for efficient clustering, noise reduction techniques are used, such as trajectory simplification.

- Evaluation Metrics:

The suggested algorithms are evaluated using a suite of thorough metrics. The silhouette score, the Davies-Bouldin index, and the internal cluster cohesion are all examples of cluster validity indices. The agreement between the algorithm-generated clusters and ground truth clusters is also measured using external validation measures, such as the adjusted Rand index and the Fowlkes-Mallows index.

- Experimental Setup:

A computer platform with industry-standard hardware is used to conduct the experiments. Languages like Python are used to implement the algorithms, with libraries like scikit-learn being utilized for optimization and HCM clustering. In order to examine how well an algorithm performs when applied to new scenarios, the researchers used a stratified sampling method that split the dataset in half.

## IV. EVALUATION AND SIMULATION

Here we discuss how to put the proposed idea into practice and evaluate it. The proposed methodology includes an optimization strategy for grouping motion trajectories. This was accomplished by optimizing the HCM clustering method using the evolutionary techniques of particle swarm and ant colony optimization. Throughout the implementation process, the MATLAB environment was utilized. The clustering procedure on the real data set containing the routes of the passenger-transporting taxis was improved in two separate circumstances using two optimization methods based on this. The primary clustering algorithm was used in both cases. The itineraries incorporate the regularly acquired GPS location data of the taxis. The results related to the evaluation criteria are recorded and inputted into the Excel software at every stage of execution. The desired charts are extracted by Excel program.

### A. Datasets

On routes that convey passengers, the suggested strategy has been tried. Data is stored using GPS coordinates that are sampled every two minutes. This is completely accurate and applies to a lot of cities in China. The data utilized pertains to the whereabouts of Beijing taxis. On March 22, 2017, during the hours of 7:50 and 7:59, the data pertains to the movements of thirty thousand taxis. We acquired the information about the origin and destination of taxis and ran additional analyses on these sites, as mentioned in reference [34]. In Fig. 5, we can see the 69514 data point pairings that include origin and destination points. With these two coordinates, you may draw out patterns of trajectories. Fig. 6 shows the results of applying the angle-based segmentation method to these points and trajectories; the resulting data set contains 13584 linear segments. Fig. 5 shows the spatial representation of the distribution of pairs of starting and ending points within a latitude and longitude range of [0.19×0.4]. Passenger traffic in

a certain location of Beijing is represented by these figures. The data in question might provide a treasure trove of information useful for traffic management and control. When constructing new streets or tearing down or fixing old ones, these numbers can also be highly useful. We also run tests of the proposed methodology on these data to see how well it works.

### B. Clustering Evaluation Criteria

Clustering algorithms can be evaluated using a variety of criteria; however, trajectory clustering requires its own set of criteria due to its distinct characteristics. Considered with other evaluation criteria such as the DB-index [14], Dunn's index [25], and XB index [26] for assessing the clustering of geographical patterns, the MBP-index [15] emerges as the most accurate and relevant. Use this metric to compare the proposed approach to other clustering techniques. The value of MBP for clustering with K classes is represented by Equation (9).

$$MBS(D) = \left(\frac{1}{D} \times \frac{F_1}{F_D} \times K_D\right)^2 \qquad (9)$$

where $K_D$ and $F_D$ are calculated from relations (10) and (11) respectively.

$$F_D = \sum_{d=1}^{D} \sum_{j=1}^{S} Hausd(C_j, s_d) \qquad (10)$$

$$K_D = max_{i,j=1}^{D}\{k(s_i, s_j)\} \qquad (11)$$

Where $k(s_i, s_j)$ is the distance between two cluster centres and $Hausd(C_j, s_d)$ is the distance between the sub-trajectory of $C_j$ and the cluster centre $c_k$. The segmentation of the data location is not crucial, and pre-specifying it is unneeded, since the MBP evaluation criterion is considered an unsupervised evaluation criterion. That is to say, the results of the clustering routes are unrelated to the road map or the actual roads.



Fig. 5. Points of origin and destinations for taxis and their locations.

Fig. 6. The use of angles for movement line segmentation.

## C. *Methods for Comparing and Assessing the Proposed Approach*

Before introducing an algorithm, it is essential to choose the correct method to evaluate and compare it. All of the method's pros and downsides should be taken into account in fair evaluations. To implement a novel approach to grouping motion routes based on optimisation of the HCM clustering algorithm, the ant colony algorithm and two-particle swarm optimisation methods were used separately in the proposed method. Consequently, the performance of each part is examined and evaluated independently. But these two suggested algorithms are tested against HCM, KMeans, and basic HCM algorithms to see how well they perform [35, 36]. An innovative approach to classifying motion trajectories, the HCM algorithm has been refined from its foundation in the cubic regression model and the Lagrange equations. In what follows, you will see the results of these comparisons according to several criteria.

## V. EVALUATION OF MBP VALUE FOR DIFFERENT K'S

When evaluating the clustering process, one of the most important parameters is the number of clusters, denoted as K. Experiments have been carried out for K values ranging from 10 to 80. What this means is that the number of clusters, as defined by the proposed algorithm and other methods, is determined before the clustering process even begins. One important part of clustering techniques is setting up the cluster centres [37]. When AGNES or HCM performs a clustering method with K clusters, the centres of these clusters are originally chosen at random, and this selection process affects the clustering results. Therefore, in order to deliver a more precise assessment of the proposed method and other results, each algorithm is executed in a 20-step process. Table I displays the results of the minimum, maximum, and average MBP values for different values of K for all methods.

TABLE I. CLUSTERING OUTCOMES COMPARED FOR VARYING K-VALUES

|  | AGNES | HCM | FCML | HCMMOPSO | HCMGA |
|---|---|---|---|---|---|
| *K=10* | | | | | |
| Max | 0.070 | 0.079 | 0.089 | 0.0770 | 0.076 |
| Mean | 0.057 | 0.079 | 0.086 | 0.0774 | 0.075 |
| Min | 0.043 | 0.078 | 0.084 | 0.0760 | 0.074 |
| *K=20* | | | | | |
| Max | 0.043 | 0.050 | 0.057 | 0.0670 | 0.059 |
| Mean | 0.037 | 0.050 | 0.057 | 0.0665 | 0.058 |
| Min | 0.030 | 0.050 | 0.056 | 0.0660 | 0.058 |
| *K=40* | | | | | |
| Max | 0.032 | 0.036 | 0.040 | 0.041 | 0.043 |
| Mean | 0.025 | 0.034 | 0.039 | 0.041 | 0.042 |
| Min | 0.019 | 0.033 | 0.038 | 0.041 | 0.052 |
| *K=80* | | | | | |
| Max | 0.025 | 0.031 | 0.034 | 0.036 | 0.039 |
| Mean | 0.022 | 0.031 | 0.033 | 0.035 | 0.039 |
| Min | 0.017 | 0.030 | 0.032 | 0.034 | 0.038 |

Fig. 7 to 10 display the overarching findings from comparing the suggested method to alternative ways for various K. By comparing the outcomes produced by the particle swarm optimisation and ant colony optimisation approaches, we can observe that the suggested method outperforms HCM, albeit to a lesser extent, for small k, i.e., when K = 10. The results show that HCMMOPSO outperforms HCMGA, the other recommended approach. However, when K is larger, the two scenarios of the suggested method outperform HCM and the two fundamental methods. Getting caught in the local minimum becomes more of an issue and the optimisation problem becomes significantly more hard as the number of clusters increases. However, the suggested approach combines the clustering algorithm with optimisation methods in an effort to achieve optimal clustering. Furthermore, when K grows, the HCMGA scenario outperforms HCMMOPSO. We will further study the ant colony optimisation approach, which outperforms the particle swarm with respect to issue complexity but has slower convergence.



Fig. 7. Evaluation of MBP values during 20 iterations with K=10.



Fig. 8. Evaluation of MBP values during 20 iterations with K=20.

**K=40**

Fig. 9.  Evaluation of MBP values during 20 iterations with K=40.

**K=80**

Fig. 10. Evaluation of MBP values during 20 iterations with K=80.

### A.  Clustering Procedure Termination Criteria

Clustering techniques iteratively update the cluster centres at each iteration to obtain the best clusters. But the most important thing is to figure out when this process should conclude. There is more than one way to finish the job. For example, it is possible to perform the same number of clustering steps if the initial number of steps is fixed. This strategy, however, appears to be somewhat illogical. On the other hand, you may consider the evaluation function and how the algorithm's efficiency is improving at each stage as an alternative approach. This value must remain over a specified threshold for the stages to be executed; otherwise, they will not be executed. The proposed method use equation (12) to halt the clustering procedure.

$$\frac{I(h+1)-I(h)}{I(h)} < \varepsilon \qquad (12)$$

In the proposed approach, the threshold value is assumed to be 0.002, and I(h) represents the value of the evaluation function at step h. The MBP standards are used to compute the I(h) value. There will be no end to the clustering processes until relation (12) is created. Nevertheless, several algorithms can achieve this with different iterations. As the algorithm gets better at meeting these requirements, its convergence speed goes up. Convergence speed alone will not be enough to achieve better performance, though. Since this rapid convergence may also lead to the local optimum, it demands consideration. As a result, optimal performance and rapid convergence are compatible. In order to compare the

convergence rates of the two situations and different techniques, we considered the execution stage of the algorithms as 100 steps and used equation (13) to get the normalised value of I(h). This criterion is used to compare the rate of convergence; after normalisation, its value will range from zero to one.

$$\frac{I - I_{min}}{I_{max} - I_{min}} \qquad (13)$$

Where are the maximum and minimum values of the evaluation function, $I_{max}$ and $I_{min}$, respectively. Fig. 11 to 14 show a comparison of the results for different cluster densities with respect to the rate of convergence. For the case where there are precisely 10 clusters, this comparison is displayed in

Fig. 14. The AGNES approach converges rapidly, as seen, but the best result it produces is inacceptably low for our issue. Steps 5 and 6 mark the end of the threshold technique, which, when applied to the question of whether the job is done, produces a locally optimal value. As a basic clustering algorithm, the HCM approach converges at the slowest rate. The HCM method, however, differs in its performance from the two alternatives. The HCM algorithm does not converge as quickly as HCMMOPSO. Almost as fast as HCMGA's convergence speed is HCM. So, it's safe to say that the particle swarm method achieves the fastest convergence when applied to clustering optimisation. With PSO's convergence speed, this conclusion was close. Nevertheless, the HCMGA method outperforms the optimal value for the clustering process.



Fig. 11. Speed of convergence comparison for K=10.



Fig. 12. Speed of convergence comparison for K=20.

Fig. 13. Speed of convergence comparison for K=40.



Fig. 14. Speed of convergence comparison for K=80.

The convergence process for 20 clusters is compared in Fig. 12, while Fig. 13 and 14 display the similar process for clusters 40 and 80, respectively. The tests also demonstrate that HCMGA converges at the fastest rate, with the convergence speeds of HCM and HCM approaches being nearly identical.

Use of a data set pertaining to the movement lines of taxis in one of China's cities, which includes the current GPS locations of taxis, allowed for the evaluation of the suggested method. Various numbers of clusters have been clustered in order to provide a more precise assessment of the suggested strategy. In situations with a small number of clusters, the findings demonstrate that the suggested method performs nearly as well as the HCM method. However, the suggested method's performance improves as the number of clusters increases. Optimisation using ant colonies achieves better clustering accuracy, while optimisation using particle swarms achieves faster convergence, according to the data.

Using HCM clustering with Ant Colony Optimisation (HCMGA) and HCM clustering with Particle Swarm Optimisation (HCMMOPSO), this study found that trajectory clustering became much more accurate and efficient. A clear and concise presentation of the results highlights the important contributions of the algorithms that were suggested.

- Clustering Accuracy Comparison:

In the first part of our investigation, we compare HCMMOPSO and HCMGA to other HCM clustering methods to see how well they cluster data. These metrics—internal cluster cohesiveness, Davies-Bouldin index, and silhouette score are presented in Table II for every algorithm. In terms of detecting significant trajectory clusters, HCMMOPSO and HCMGA both routinely beat baseline approaches.

TABLE II.        CLUSTERING ACCURACY METRICS

| Metric | HCMMOPSO | HCMGA | Baseline 1 | Baseline 2 |
|---|---|---|---|---|
| Silhouette Score | 0.71 | 0.80 | 0.61 | 0.53 |
| Davies-Bouldin Index | 0.40 | 0.31 | 0.54 | 0.60 |
| Internal Cohesion (avg) | 0.91 | 0.92 | 0.83 | 0.80 |

Higher silhouette scores and lower Davies-Bouldin indices show that the clusters are better defined and well-separated, indicating a significant improvement in clustering quality. HCMMOPSO and HCMGA demonstrate exceptional internal cohesiveness, which is a direct result of the algorithms' ability to produce dense clusters.

- Statistical Significance:

With p-values lower than the 0.05 threshold, statistical analyses, such as ANOVA and t-tests, validate the importance of the noticed improvements. Statistics confirm that HCMMOPSO and HCMGA are better than baseline approaches, which gives the suggested trajectory clustering method more credibility.

- Comparison to Existing Studies:

We compare our results to those of related studies so that you can put them in context. Importantly, as compared to studies that used conventional clustering techniques, our silhouette scores are far higher. The optimisation procedures that were introduced allow for more accurate trajectory segmentation, which in turn improves the performance of HCMMOPSO and HCMGA.

- Discussion:

Results showing an increase in clustering accuracy show that HCMMOPSO and HCMGA are useful in the real world, especially for optimising taxi dispatch systems. The overall efficiency of transport services is improved by effectively identifying and differentiating between distinct trajectory patterns.

Finally, HCMGA and HCMMOPSO are great examples of how trajectory analysis has progressed thanks to the integration of HCM clustering and optimisation approaches. This study's findings add to the expanding literature on trajectory clustering and point the way towards further investigations on transportation system optimisation. Our approach has the potential to improve the efficiency of trajectory-based applications, and the shown improvements highlight its practical value.

## VI. CONCLUSION

An approach based on the HCM clustering algorithm for route clustering was proposed as a means to extract city passenger flow patterns. Using a sequence of GPS markers that the user specifies, this method plots a route that is accessible to both pedestrians and drivers. The data shown here pertain to a journey that starts in one place and finishes in another. At regular intervals (often less than two minutes) during this process, the location of the moving object is captured using sampling. This is how the information about the moving cab's many points is recorded. The suggested method specifies the route information as a pair consisting of the geographical dimensions of the spot and the angle of the line that connects it to the next place. Before extracting the sublines, the data undergo pre-processing as part of the clustering process. An initial step in improving clustering accuracy and speed is to partition the sublines according to the angle between them. After that, HCM is used to cluster the data. Cluster centres were found in the local optimum, and the HCM algorithm failed to optimise the objective function, according to the evaluation and implementation of the algorithm. In order to fix this, we talked about meta-heuristic algorithm research and how to apply it to improve the HCM algorithm.

In two different cases, the HCM clustering method was fine-tuned using ant colony and particle swarm methods. Using a real data set containing cab movement data in Beijing, China, and the MATLAB environment, the suggested method was implemented. To test and compare the suggested method, a number of experiments were conducted. By applying the MBP criterion to the clustering accuracy performance, we found that the two scenarios outperformed other methodologies as the number of clusters increased. For a large number of clusters, the evaluation results show that the proposed method outperforms existing approaches in terms of clustering accuracy. The HCMGA algorithm improved clustering accuracy for MBP values by 0.61% compared to prior approaches, while the HCMMOPSO methodology improved it by 0.64%. The findings of the implementation also show that the particle swarm method converges faster and the ant colony approach is the most accurate.

A specific dataset concerning cab traffic in a Chinese city, namely Beijing, is used to test the proposed method. Data collected from GPS coordinates every two minutes form the basis of the evaluation. As an additional metric for cluster evaluation, the MBP index is utilised in the assessment.

For future studies, the proposed method might be expanded and generalised in other ways. For example, this method can be employed to classify various motion trajectories. Alternately, use more evolutionary optimisation methods to improve grouping. You may also test how well the method works by changing the pre-processing and sub-line segmentation steps.

## REFERENCES

[1] H.-L. Ling, J.-S. Wu, Y. Zhou, and W.-S. Zheng, "How many clusters? A robust PSO-based local density model," Neurocomputing, vol. 207, pp. 264–275, 2016.

[2] Y. Wu, "A survey on population-based meta-heuristic algorithms for motion planning of aircraft," Swarm Evol Comput, vol. 62, p. 100844, 2021.

[3]     M. S. Choudhry and R. Kapoor, "Performance analysis of fuzzy C-means clustering methods for MRI image segmentation," Procedia Comput Sci, vol. 89, pp. 749–758, 2016.

[4]     Yumeng Cao, Ning Xu, Huanqing Wang, Xudong Zhao, A. Ahamed. Neural Networks-Based Adaptive Tracking Control for Full-State Constrained Switched Nonlinear Systems With Periodic Disturbances and Actuator Saturation, International Journal of Systems Science, 54(14): 2689-2704, 2023.

[5]     Y. R. Alekya Rani and E. Sreenivasa Reddy, "An optimal communication in WSN enabled by HCM clustering and improved meta-heuristic model," International Journal of Pervasive Computing and Communications, vol. 19, no. 2, pp. 233–254, 2023.

[6]     Pedditi, R. B., & Debasis, K. (2024). MACR: A Novel Meta-Heuristic Approach to Optimize Clustering and Routing in IoT-based WSN. International Journal of Intelligent Systems and Applications in Engineering, 12(1), 346-359.

[7]     Tengda Wang, Liang Zhang, Ning Xu, and K. H. Alharbi. Adaptive critic learning for approximate optimal event-triggered tracking control of nonlinear systems with prescribed performances, International Journal of Control, https://doi.org/10.1080/00207179.2023.2250880, 2023.

[8]     Wang, G., Wu, J., & Trik, M. (2023). A Novel Approach to Reduce Video Traffic Based on Understanding User Demand and D2D Communication in 5G Networks. IETE Journal of Research, 1-17.

[9]     J. Sun, Y. Zhang, and M. Trik, "PBPHS: a profile-based predictive handover strategy for 5G networks," Cybern Syst, pp. 1–22, 2022.

[10]    M. Trik, H. Akhavan, A. M. Bidgoli, A. M. N. G. Molk, H. Vashani, and S. P. Mozaffari, "A new adaptive selection strategy for reducing latency in networks on chip," Integration, vol. 89, pp. 9–24, 2023.

[11]    Ding, X., Yao, R., & Khezri, E. (2023). An efficient algorithm for optimal route node sensing in smart tourism Urban traffic based on priority constraints. Wireless Networks, 1-18.

[12]    Yue, S., Niu, B., Wang, H., Zhang, L. and Ahmad, A.M. (2023), "Hierarchical sliding mode-based adaptive fuzzy control for uncertain switched under-actuated nonlinear systems with input saturation and dead-zone", Robotic Intelligence and Automation, Vol. 43 No. 5, pp. 523-536. https://doi.org/10.1108/RIA-04-2023-0056.

[13]    Kadir, D. H. (2021). Statistical evaluation of main extraction parameters in twenty plant extracts for obtaining their optimum total phenolic content and its relation to antioxidant and antibacterial activities. Food Science & Nutrition, 9(7), 3491-3499.

[14]    Xiao, L., Cao, Y., Gai, Y., Khezri, E., Liu, J., & Yang, M. (2023). Recognizing sports activities from video frames using deformable convolution and adaptive multiscale features. Journal of Cloud Computing, 12(1), 1-20.

[15]    Omer, A. W., Blbas, H. T., & Kadir, D. H. (2021). A Comparison between Brown's and Holt's Double Exponential Smoothing for Forecasting Applied Generation Electrical Energies in Kurdistan Region.

[16]    M. Braik, M. H. Ryalat, and H. Al-Zoubi, "A novel meta-heuristic algorithm for solving numerical optimization problems: Ali Baba and the forty thieves," Neural Comput Appl, vol. 34, no. 1, pp. 409–455, 2022.

[17]    M. Trik, A. M. N. G. Molk, F. Ghasemi, and P. Pouryeganeh, "A hybrid selection strategy based on traffic analysis for improving performance in networks on chip," J Sens, vol. 2022, 2022.

[18]    Saleh, D. M., Kadir, D. H., & Jamil, D. I. (2023). A Comparison between Some Penalized Methods for Estimating Parameters: Simulation Study. QALAAI ZANIST JOURNAL, 8(1), 1122-1134.

[19]    J. Riedy, "Updating pagerank for streaming graphs," in 2016 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), IEEE, 2016, pp. 877–884.

[20]    Shen, Y. (2024). Robotic trajectory tracking control system based on fuzzy neural network. Measurement: Sensors, 101006.

[21]    Sajadi, S. M., Kadir, D. H., Balaky, S. M., & Perot, E. M. (2021). An Eco-friendly nanocatalyst for removal of some poisonous environmental pollutions and statistically evaluation of its performance. Surfaces and Interfaces, 23, 100908.

[22]    W. Qiao and Z. Yang, "An improved dolphin swarm algorithm based on Kernel Fuzzy C-means in the application of solving the optimal problems of large-scale function," Ieee Access, vol. 8, pp. 2073–2089, 2019.

[23]    Blbas, H., & Kadir, D. H. (2019). An application of factor analysis to identify the most effective reasons that university students hate to read books. International Journal of Innovation, Creativity and Change, 6(2), 251-265.

[24]    Chen Cao, Jianhua Wang, Devin Kwok, Zilong Zhang, Feifei Cui, Da Zhao, Mulin Jun Li, Quan Zou. webTWAS: a resource for disease candidate susceptibility genes identified by transcriptome-wide association study. Nucleic Acids Research.2022, 50(D1): D1123-D1130.

[25]    Fakhri, P. S., Asghari, O., Sarspy, S., Marand, M. B., Moshaver, P., & Trik, M. (2023). A fuzzy decision-making system for video tracking with multiple objects in non-stationary conditions. Heliyon.

[26]    Y. Zheng, "Trajectory data mining: an overview," ACM Transactions on Intelligent Systems and Technology (TIST), vol. 6, no. 3, pp. 1–41, 2015.

[27]    Kadir, D. (2018). Bayesian inference of autoregressive models (Doctoral dissertation, University of Sheffield).

[28]    Haoyu Zhang, Quan Zou, Ying Ju, Chenggang Song, Dong Chen. Distance-based Support Vector Machine to Predict DNA N6-methyladine Modification. Current Bioinformatics. 2022, 17(5): 473-482.

[29]    Wang, G., Wu, J., & Trik, M. (2023). A Novel Approach to Reduce Video Traffic Based on Understanding User Demand and D2D Communication in 5G Networks. IETE Journal of Research, 1-17.

[30]    Khezri, E., Yahya, R. O., Hassanzadeh, H., Mohaidat, M., Ahmadi, S., & Trik, M. (2024). DLJSF: Data-Locality Aware Job Scheduling IoT tasks in fog-cloud computing environments. Results in Engineering, 101780.

[31]    Heng Zhao, Huanqing Wang, Ning Xu, Xudong Zhao, Sanaa Sharaf, Fuzzy approximation-based optimal consensus control for nonlinear multiagent systems via adaptive dynamic programming, Neurocomputing, 533: 126529, 2023.

[32]    Trik, M., Boukani, B., Ansari, B., Emtiyaz, S., Azar, S. R., & Mohammadi, F. (2014). An Overview of through-silicon via–based three dimensional integrated circuits (3D IC) to placement to optimize timing. Research Journal of Recent Sciences.

[33]    Khezri, E., Zeinali, E., & Sargolzaey, H. (2023). SGHRP: Secure Greedy Highway Routing Protocol with authentication and increased privacy in vehicular ad hoc networks. Plos one, 18(4), e0282031.

[34]    Khosravi, M., Trik, M., & Ansari, A. (2024). Diagnosis and classification of disturbances in the power distribution network by phasor measurement unit based on fuzzy intelligent system. The Journal of Engineering, 2024(1), e12322.

[35]    Pedroche, D. S., Herrero, J. G., & López, J. M. M. (2024). Context learning from a ship trajectory cluster for anomaly detection. Neurocomputing, 563, 126920.

[36]    Liu, L., Wang, X., Wang, X., Xie, J., Liu, H., Li, J., ... & Yang, X. (2024). Path Planning and Tracking Control of Tracked Agricultural Machinery Based on Improved A* and Fuzzy Control. Electronics, 13(1), 188.

[37]    Kapnopoulos, A., Kazakidis, C., & Alexandridis, A. (2024). Quadrotor trajectory tracking based on backstepping control and radial basis function neural networks. Results in Control and Optimization, 14, 100335.

# Sustainability and Resilience Analysis in Supply Chain Considering Pricing Policies and Government Economic Measures

Dounia SAIDI[1], Aziz AIT BASSOU[2], Jamila EL ALAMI[3], Mustapha HLYAL[4]

LASTIMI Laboratory, High School of Technology in Sale, Mohammed V University in Rabat, Rabat, Morocco[1, 2, 3]
Logistics Center of Excellence, Higher School of Textile and Clothing Industries, Casablanca, Morocco[4]

*Abstract*—Sustainability and resilience are becoming increasingly critical in shaping supply chain pricing strategies. They ensure that supply chains can withstand disruptions while adhering to environmental and social standards, thereby securing long-term economic viability. Despite their importance, the integration of these two pillars with the promotion of domestic products remains under-explored, especially concerning their influence on the competitive dynamics within supply chains. This study seeks to bridge this gap by examining the influence of sustainability, resilience, and domestic product promotion on supply chain pricing strategies. We introduce a model that captures the interactions among a central supplier, multiple stores, and the government, focusing on strategies adopted by each stakeholder to maximize its profit while adhering to sustainability and resilience requirements. The study reveals that stores' pricing strategies are significantly influenced by their sustainability efforts, with the cost coefficient of these efforts and the elasticity of sustainability efforts directly affecting profit margins. It also finds that the supplier's resilience strategy involves allocating inventory reserves to manage wholesale pricing effectively. Governmental regulatory measures, through taxation and subsidies, are shown to play a crucial role in maintaining the balance between domestic and foreign products and providing flexibility to diversify product sources to cope with local disruptions. Finally, perspectives are provided to enrich the understanding of how sustainability and resilience can be considered and impact pricing policies of the whole network.

*Keywords—Supply chain management; pricing policies; sustainability; resilience; government regulation*

## I. INTRODUCTION

In today's interconnected world, supply chains stand as the backbone of global commerce, ensuring the seamless flow of goods and services across continents. However, as we navigate the challenges of the 21st century, from climate change and resource scarcity to geopolitical tensions and technological disruptions, the importance of sustainable and resilient supply chains has never been more pronounced [1]. Sustainability in supply chains ensures that operations are conducted in an environmentally conscious, socially responsible, and economically viable manner. On the other hand, resilience equips supply chains with the agility and adaptability to withstand unforeseen challenges, be they natural disasters, trade restrictions, or global pandemics. A supply chain that embodies both these qualities not only ensures business continuity and profitability but also plays a pivotal role in fostering a sustainable future for all.

This study embarks on a meticulous exploration of these twin pillars—sustainability and resilience—within the context of a supply chain of one supplier and several stores. It underscores the challenges tied to bolstering local production in an era dominated by transnational logistics networks. In today's interconnected world, gaining a nuanced understanding of how supply chains can optimize profitability while simultaneously fostering positive local and environmental impacts is imperative. Thus, the dual role of a supply chain—as a catalyst for economic gains and as an environmentally and locally attuned entity—warrants in-depth scrutiny [2], [3]. Firstly, it is a matter of drawing out the possible links between sustainability and resilience within the supply chain [3]. As such, the relationship between sustainability initiatives and the supply chain's ability to cope with disruption needs to be specifically addressed, exploring how sustainable practices can strengthen resilience. Next, it is necessary to know how a store can maximize profits while maintaining sustainable practices in its supply chain. This can encompass pricing strategies, operational efficiencies, and the integration of sustainable practices to ensure profitability [4]. Finally, the focus is on approaches to actively favor local products while meeting sustainability and resilience requirements. The main objective is to explore the measures needed to encourage local manufacturers, manage demand, and establish cooperation between the various players in the supply chain.

The interest of this research lies in several crucial aspects. Firstly, it is important to note that the issue of promoting domestic products while maintaining sustainability and resilience in supply chains remains relatively underexplored in the scientific literature. Very few studies to date have addressed this complex issue, despite its growing importance in a world where the globalization of supply chains is increasingly being called into question.

In addition, the recent crisis of supply chain disruption, exacerbated by unpredictable events such as the COVID-19 pandemic, has highlighted the urgency of rethinking and strengthening supply chain resilience [5]. This crisis has also raised key questions about the vulnerability of global supply chains and the importance of promoting local production to reduce this vulnerability [6].

Consequently, this research fills a gap in the literature by addressing these interrelated issues of sustainability, resilience, and the promotion of domestic products in supply chains. It thus offers innovative perspectives to address current and future challenges faced by companies and governments in supply chain management, helping to create more robust and environmentally friendly solutions in an ever-changing world. Practical implications of the findings are significant, as they empower stakeholders such as stores and suppliers to strategically choose pricing strategies that not only maximize profit but also promote sustainability and domestic products, enhancing resilience in the face of disruptive events. This study provides valuable insights for companies across various sectors, including textiles, mass-market retail, automotive, and more. The ability to navigate competitive dynamics, even within monopolistic scenarios, is crucial, especially when a company controls a substantial market share. In instances where a manufacturer (the monopolist) also owns distribution channels or outlets, competition with independent stores becomes a reality. Understanding this complexity, the study aims to provide answers to the following research questions.

Given this complexity, this study aims to answer the following research questions.

- What are pricing strategies to adopt by the store and the supplier to maximize their profits considering sustainability?

- How can the supplier comply with sustainability and resilience requirements?

- How can the government promote sustainability by encourage domestic products and resilience by offering the possibility of importing products?

The present research work will focus on a sustainable and resilient two-tier monopoly supply chain framework. Our supply chain management model examines the interactions between a central supplier, several stores and the government. The stores determine pricing and sustainability strategies, the supplier manages distribution and reserves, while the government regulates via taxes and subsidies, particularly with regard to domestic and foreign products. The key element is the resilience of the chain, with considerations such as inventories, the social burden of stores and the diversification of supply sources. The aim is to analyze supply chain dynamics with a focus on sustainability and local production.

The paper is structured as follows: Section II explores the literature related to supply chain pricing strategies, with a focus on promoting sustainability, resilience and domestic products with a government intervention. Section III concerns the model construction and analysis with three pricing strategies, the first one concern the store level, the second one concern the supplier level and finally the government. Section IV conducts a numerical analysis using examples to support the choice of pricing strategies adopted at each level. Moving to Section V, we present and meticulously discuss the findings derived from the developed model and the numerical analysis. Section VI encapsulates these findings in a comprehensive conclusion, elucidates the study's limitations and offers final observations and prospective directions.

## II. LITERATURE REVIEW

Academic research is increasingly exploring these intersections between sustainability, resilience, performance, and competitiveness, seeking to define the ways in which companies can maintain efficient operations while being ecologically responsible and resilient to disruption.

This work is closely related to sustainability and resilience, domestic and foreign product and governmental intervention pricing strategies for supply chains.

### A. Sustainability Pricing Strategies

The first focus of this paper is on the articulation and implementation of sustainability strategies. Firms are channeling resources into social and environmental programs to align with sustainability mandates, driven by a mix of governmental directives reflecting a heightened awareness [7]–[10]. In the study [11] authors investigated how inventory restocking and sustainability funding are influenced by three distinct regulatory environments. Similarly, the work in [12] investigated the best production practices for various products within the framework of cap-and-trade regulations, and in study [13] honed in on carbon emissions from warehouses, analyzing the inventory management and investment in eco-friendlier technologies in response to carbon emission limits set by cap-and-trade policies. While the study. Within the scope of this research, it's commonly acknowledged that environmental taxes significantly motivate corporate investment towards sustainable practices [14].

On the other hand, investing in sustainability initiatives is a strategic marketing proposition to enhance the company's brand image, examines various factors that impact sustainability of supply chain. Consumers' environmental awareness, governmental regulations, sustainability investments, pricing, production quantities, and environmental constraints interact and influence each other within the context of sustainability and its impact on products and production processes [15]. Likewise, the investigation of retailers' investments in environmental R&D and manufacturers' proposed strategies for balancing and coordinating the supply chain through various joint R&D contracts. The study also reveals that consumer environmental awareness, while not always leading to increased demand for green products, systematically boosts the profits of green supply chains [16]. In the study [17] authors explores the impact of cooperative promotion on decisions and sustainability of service platform supply chains from different markets. Theoretical models suggest that joint promotion is advantageous, especially when independent promotional activities moderately impact demand. However, the benefit of cooperative promotion, influenced by demand sensitivity and adjustments in price and quality, may vary. Platforms with a high baseline demand typically derive more benefits from cooperation but might be less willing to invest more in cooperative promotions.

The examination of the relationship between sustainability efforts or investments and various aspects of product demand, production, or supply chain management, were explored in [18] where the manufacturer's incentive is considered to reduce carbon emissions in the presence of carbon taxes under revenue-sharing and cost-sharing contracts. Meanwhile, a

comparison of optimal pricing and sustainability efforts was conducted in two scenarios in the study [19]: one where the manufacturer is a for-profit company driven by profit-seeking motives, and the other where the manufacturer is a nonprofit organization aiming to maximize demand realization and quantity. The same approach was made by [11], which examines sustainability investments made by retailers. Differing from this line of research, our paper considers both scenarios, where sustainability investments can be initiated by either the manufacturer or the retailer. Similarly, optimal decisions in pricing and ecological investment within competitive supply chains for electric vehicles were explored, focusing on the tension between the costs and revenues of green technologies. Manufacturers and retailers, in various market scenarios modeled via game theory approaches, navigate between investments in ecological technologies and consumer sensitivity to prices and ecology. The results offer insights to guide electric vehicle companies in developing optimal green investment and pricing strategies [20].

### B. Resilience Pricing Strategies

The literature on resilience strategies in supply chain pricing is a dynamic and evolving field that addresses the crucial need for businesses to adapt and thrive in an increasingly uncertain and complex environment. Researchers and practitioners alike have recognized that pricing decisions are not just about setting optimal prices but also about fortifying supply chains against disruptions and unforeseen challenges.

Many studies underscore the critical importance of resilience in supply chain management and the diverse approaches to achieving it in the face of environmental, operational, and market challenges. A complex supply chain involving three manufacturers and a distributor managing complementary and substitute products which emphasizes the resilience of the chain against various possible disruptions. Using game theories to determine optimal prices at different levels of the supply chain, the developed model seeks to navigate efficiently through these potential interruptions, thereby ensuring order fulfillment and system stability despite environmental and operational challenges [21]. As well as, [22] where another strategy based on maintaining extra inventory at distribution centers is implemented and ensuring the reliability of distribution centers, which positively impacted the competitiveness and adaptability of supply chains investigated the influence of resilience strategies in supply chain management, particularly within the context of price competition and facility disruptions. In another context related to insufficient capacity, authors of the work [23] scrutinizes resilience in the maritime supply chain, specifically focusing on the "co-opetition" relationship between shipping companies and freight forwarders. Using a model based on game theory, it reveals that establishing a direct sales platform by shipping companies strengthens their competitive position and improves their price and profit compared to freight forwarders. Moreover, in the event of capacity shortage, the strategic implementation of capacity allocation and pricing strategies, especially through a spot market, can enhance the resilience of the supply chain. The examination of resilient agricultural supply chains in the post-COVID-19 era, as explored in [6],

delves into the utilization of channel leadership strategies. This study emphasizes the critical need for selecting appropriate leadership tactics to ensure maximum profitability, optimal pricing, and high service quality in a volatile market. It also offers practical insights for the transition to e-commerce platforms in the aftermath of the pandemic.

By embracing both foreign and domestic products, businesses can diversify their portfolios, ensuring a balanced approach to sourcing and mitigating risks associated with supply disruptions. Studies in the literature delve into how these models can be leveraged to not only bolster the profitability of domestic enterprises but also strengthen their position in the increasingly competitive and complex market landscape. The optimization of prices and profits for a domestic company, along with the reduction of retailer costs, is the focus of the research [24]. By examining various scenarios, the research establishes that implementing an adapted pricing strategy can significantly enhance the profitability and market competitiveness of the domestic company in relation to imported products. Similarly, the pricing competition between national and foreign manufacturers on diversified market segments, using a Stackelberg game model was explored by [25] taking into account factors such as price and quality to influence customer purchasing trends, the study concludes that market segmentation by income levels can increase profits for the national manufacturer and improve its competitive advantage against the foreign manufacturer. Likewise, a cooperative strategy between national pharmaceutical manufacturers and foreign licensors to exclusively produce locally licensed branded drugs, with the aim of increasing the manufacturer's revenues and potentially offering government discounts [26]. Using a Nash bargaining solution and cooperative game theory to model tariff negotiations, the research reveals that the local market share of the licensor and the return on capital of national manufacturers positively influence the equilibrium price. If the price is too low, foreign entities might lack the incentive to join the coalition. On the other hand, in the context of supply disruption risk the study [27] examines the competitive dynamics between two closed supply chains, focusing on the management of product prices and recycling. One of the supply chains, a retailer, has the choice between a reliable but expensive domestic supplier and a cheaper but unreliable foreign supplier. The results indicate that in a dual supply situation, there is a direct correlation between sourcing from a foreign supplier and the return rate of used products, and that strategic use of return policies is essential to maintain competitiveness in the market.

### C. Government Regulation

The government intervention strategies in various sectors are mainly focusing on subsidies, taxes, and regulations to promote environmentally friendly practices and sustainable development. To stimulate remanufacturing activities, governmental subsidies was explored, revealing that excessively high or low financial aids prompt remanufacturers to compete with producers [28]. A case study involving five European countries was highlighted by [29], unveiling how the evolution of the green economy is shaped by governmental intervention. Similarly, the price competition between green and non-green products is analyzed including government

subsidies and the implementation of taxes to minimize carbon emissions [30]. Whereas, the influence of the government subsidy on the green supply chain system, exploring various chain leadership modes and how the subsidy policy impacts each chain member and the system's profit [31]. Meanwhile, the chain structure and pricing decisions for the producer and government subsidy strategy is studied, contrasting new and remanufactured products [32]. The impact of government incentive strategies, approaching eco-responsible products from a game theory-based perspective and comparing the benefits of each chain member, the level of eco-friendliness, and environmental improvement [9]. The assessment of the impact of governmental interventions on bioenergy and conventional energy supply chains is studied in [33], revealing that some support strategies, especially investment subsidies, can significantly optimize both profits and carbon emission reduction efforts while supporting sustainable development goals. In the research [7], authors utilized game theoretical models to refine pricing for energy sectors, aligning with government, societal, and ecological objectives. Findings reveal Nash strategies boost governmental and societal benefits, whereas cooperative approaches favor ecological results and energy producers' earnings. While the government regulation strategies are explored for promoting EV adoption and reducing $CO_2$ emissions through targeted tax reforms and subsidies, indicating that government policy adjustments are essential for achieving sustainability and market influence [34].

## III. MODELING FORMULATION AND ANALYSIS

### A. Model Description and Assumptions

In our supply chain model, we consider a configuration, as illustrated in Fig. 1, where the actors in this system are the stores, the supplier and the government. Each of these actors has an impact on the supply chain through specific strategies. The central supplier provides multiple stores with foreign and domestic products. The stores, as the final sale points, make significant strategic decisions regarding prices, sustainability efforts, and demand management. The supplier, although not directly involved in the production of products, plays a central role in distribution, negotiating unit prices and managing security stock $(\psi_i)$ with each store. The government steps in by imposing custom fees $(\tau)$ on foreign products and by granting subsidies $(\nu)$ to promote domestic products as presented in the Table I. It also plays a regulatory role by influencing tariff policies and sustainability practices.

Integrating sustainability and resilience into supply chains combines environmental responsibility, operational robustness, and market competitiveness. This research paper primarily focuses on sustainability strategies, which are incorporated into our model at various levels. At the store level, sustainability efforts ( $s_i$ ) involve demanding durable products and undertaking sustainability actions and investments, such as sustainable product refurbishment, the requirement for sustainable packaging, recycling, and more [35]. At the supplier level, it includes contributing to the reduction of carbon footprint through the integration of $CO_2$ emissions cost ( $\eta_i$ ) during product transportation. Additionally, at the government level, the social charge ( $\lambda$ ) paid by the store

contributes to fostering a more socially sustainable business environment [35] across the supply chain.

Moreover, the proposed model takes into account the resilience of the supply chain, through the integration of parameters that attempt to express said resilience, such as the security stock and the diversification of supply sources. Indeed, the supplier reinforces supply chain resilience through inventories that play a key role by implementing a dedicated extra inventory $(\psi_i)$ for each store i, that is essential to deal with potential disruptions [22]. Additionally, diversifying supply sources is also a crucial aspect of resilience. The supplier can choose to diversify the supply chain, whether they are national or foreign. On one hand, importing products allows flexibility and reduce risks of supply chain disruption and increase their resilience in the face of uncertainty. On the other hand, encouraging domestic products reduce dependence on international markets, secure jobs and bolster economic security within a country. This paradox stems from the conflicting goals of achieving supply chain flexibility and resilience by depending on imports, and simultaneously encouraging local manufacturing.

The government plays a critical role in navigating this paradox through various regulatory mechanisms. In this sense, government by offering subsidies $(\nu)$ to actively promotes domestic products $(1-\theta)$, while imposing custom fees $(\tau)$ to foreign products $(\theta)$ can them more expensive. Furthermore, the government maintains the option to import foreign products $(\theta)$. This dual approach not only provides flexibility but also serves as a valuable contingency plan in the face of disruptive events, significantly enhancing the supply chain's ability to endure and recover from disruptions. These two key options are considered, each with clear financial implications:

*Import Foreign Products:* The supplier is considering importing a large volume of foreign products, such as manufactured goods or raw materials, for resale on the local market. However, this decision requires careful cost management. For example, when importing these products, you need to take into account customs costs, which vary according to the type of product and the country of origin. Customs costs can represent a significant proportion of import expenditure. A relevant example is Canada, where a survey of 635 men and women revealed positive attitudes towards Canadian-made products, particularly among women. However, it is important to note that customs costs can vary from country to country. Therefore, the supplier must calculate these costs accurately to assess the economic viability of this option.

*Purchase of domestic products:* In addition to importing, the supplier plans to purchase local products. A concrete example might be the purchase of locally manufactured products for resale on the domestic market. The government encourages this approach by offering economic benefits, such as subsidies to support local businesses or tax breaks for domestic products. These financial incentives support the purchase of local products and can reduce procurement costs. Suppliers must therefore integrate these advantages into their local purchasing strategy to maximize the economic benefits.

Fig. 1. Theoretical model of a network of stores and one supplier with a government regulation.

Profits within the supply chain are determined by the pricing and sustainability strategies set up by each store. Each store must decide the selling price of its products, taking into account various factors such as operational costs ( $c_i$ ), negotiated purchase prices ($\omega_i$) with the supplier, sustainability efforts ($s_i$) deployed, and other parameters. The government intervenes by imposing taxes on $CO_2$ emissions related to product transportation, while subsidies may be granted to encourage sustainability.

In this model, each store i negotiates the purchase price ($\omega_i$) individually with the supplier, allowing for significant flexibility and customization in business relationships. This freedom to negotiate creates a competitive situation where the supplier and each store i strategically compete to achieve an optimal purchase price ( $\omega_i$ ), considering the competing interests of each actor within the supply chain.

This supply chain model aims to understand interactions and dynamics between actors, configurations, profits, and resilience in a context of promoting local production and sustainable concerns.

*1) Assumptions:* Within the framework of our supply chain pricing model, several fundamental assumptions are established to simplify and define the context of our analysis. These assumptions define the parameters, relationships and basic conditions governing our system. They are essential for framing our study precisely and rigorously:

- $a, b, k, \alpha > 0$
- $p_i > (c_i + \omega_i)$
- $\frac{k^2}{2\alpha} > b$
- We consider a monopolistic market configuration at store level;

- A store i is supplied by a single supplier;
- We assume uniform operator numbers for all stores, as we also assume they employ the same value of $d_m$, So $\forall\, i, A_i = A$ ;
- The supplier applies different selling prices negotiated with each store i;
- The supplier keeps a reserve quantity ($\psi_i$ ) for each store i;
- Each store has fixed costs that are independent of demand (wages, rent, electricity bills, ......);
- The Supplier transports the products to the stores. This transport generates $CO_2$ emissions, which are taxed by the government;
- It is assumed that all planned requests have been met. In this case, we won't deal with profit expectations, and the model is considered deterministic;
- Products are assumed to have the same quality preference for the customer.

However, in the realm of supply chain and inventory management, the supplier's decision to reserve a specific quantity of goods $\psi_i$ for each order from a store i reflects a strategic approach to ensure the supplier's resilience in the face of market uncertainties. This quantity acts as a buffer, allowing the supplier to respond efficiently to variations in demand and unforeseen disruptions in the supply chain.

Based on the above assumptions, the objective functions of the problem to be modeled are as follows:

- Maximize store profit
- Maximize supplier profit
- Maximize government revenues according to tax/subsidy policies

Table I, presented below, succinctly encapsulates the key sustainability, resilience, and economic parameters employed in the formulated model, serving as a valuable aid for comprehension.

TABLE I. SUSTAINABILITY, RESILIENCE AND ECONOMIC PARAMETERS OF THE MODEL

|  | Sustainability | Resilience | Economic |
|---|---|---|---|
| Store | $s_i, k, \alpha, A_i$ | $\theta, \psi_i$ | $p_i, \pi_i$ |
| Supplier | $\eta_i$ | $\theta, \psi_i, \gamma$ | $\omega_i, \pi_{i,s}, \pi_s$ |
| Government | $\lambda$ | $\theta, \psi_i$ | $\tau, \nu, \pi_g$ |

*2) Model* parameters and variables

$d_i$ : demand quantity at store i
$p_i$ : unit price at store i
$s_i$ : sustainability efforts for store i
$a$ : basic market demand
$b$ : price elasticity
$k$ : elasticity of sustainability effort

$c_i$ : operating unit cost of store i

$c$ : operating unit cost of supplier

$F_i$ : fixed expenses of store i

$\alpha$ : cost coefficient of sustainability effort for store i

$d_m$ : average demand estimated by store i

$A_i$ : number of operators needed to satisfy $d_m$

$\lambda$: social charge per operator to be remitted to the government by a store

$y_i$ : unit cost to be paid to the operator to satisfy $d_m$

$\gamma$ : inventory security sensitivity

$\psi_i$ : inventory quantity reserved for store i

$\eta_i$ : unit cost of $CO_2$ emission

$\theta$ : ratio of quantity of foreign products, $\theta \epsilon [0,1]$

$\tau$ : custom fees

$v$ : government subsidy

$\omega_i$ : wholesale price negotiated between supplier and store i

$w_i^s$ : wholesale price negotiated considering the whole network

$\pi_i$ : profit of store i

$\pi_{i,s}$ : profit of the supplier considering the whole network

$\pi_s$ : total supplier profit

$\pi_g$ : total government profit

### B. Model Construction and Analysis

In economic analysis, the inverse demand function emerges as a pivotal tool, offering a distinct perspective compared to the conventional demand function, which typically represents quantity demanded as a function of price [36]. This inverse approach articulates the price as contingent upon the quantity demanded, essentially inverting the traditional relationship. For companies, this analytical approach is indispensable for formulating sophisticated pricing strategies.

Similar to the work of [37], [38] the inverse demand equation is given as follows:

$$d_i = a - bp_i + ks_i \qquad (1)$$

Based on the work of [39], [40] $C(s_i) = \frac{\alpha s_i^2}{2}$, corresponds to the cost of the sustainability effort.

*1) Store i strategies:* In this sub-section, the store's profitability is studied. So, it can be divided into two components: the first part comprises gains, represented by the product of the price $(p_i)$ and demand $(d_i)$, while the second part encompasses expenses. Among these expenses, some are variable, such as operational costs $(c_i)$ and the negotiated wholesale price between the supplier and store i $(\omega_i)$, while others are fixed, such as overhead expenses $(F_i)$, the employee payroll $(A_i y_i)$, and sustainability efforts $(s_i)$.

The profit Eq. (2) for store i is given as follows.

$$\pi_i = d_i(-c_i + p_i - \omega_i) - F_i - \frac{\alpha s_i^2}{2} - A_i y_i \qquad (2)$$

By replacing Eq. (1) in Eq. (2) we get:

$$\pi_i(p_i, s_i) = -bp_i^2 + p_i(a + ks_i + b\omega_i) - \frac{\alpha s_i^2}{2} -$$

$$c_i(a - bp_i + ks_i) - (a + ks_i)\omega_i - F_i - A_i y_i \qquad (3)$$

According to this equation, a store has two strategies: The price to apply and the sustainability effort to adopt.

Assuming now that the store focuses only on the price strategy. The price to achieve maximum profit is:

$$p_i^* = \frac{1}{2}\left[\frac{a+ks_i}{b} + (c_i + \omega_i)\right] \qquad (4)$$

It can be seen that the selling price increases as the sustainability effort increases. It is also influenced by the purchase price $\omega_i$ and operating costs.

Price elasticity tends to eliminate the effect of sustainability effort. For this reason, we assume that $k > b$.

**Proof.** Maximum profit is sought by deriving the profit function $\pi_i$.

$$\frac{d\pi_i}{dp_i} = -2bp_i + a + b(c_i + \omega_i) + ks_i$$

We have: $\frac{d^2\pi_i}{d^2p_i} = -2b \leq 0$ thus, the function admits a maximum. This maximum is obtained by solving $\frac{d\pi_i}{dp_i} = 0$

Since $p_i^*(s_i)$ has the equation of a straight line and $\frac{dp_i^*}{ds_i} = \frac{k}{2b} > 0$. Then the function is increasing with respect to $s_i$.

Furthermore, assuming that the store wishes to use the sustainability strategy to attract more demand. The sustainability effort to achieve maximum profit is:

$$s_i^* = \frac{k(p_i - (c_i + \omega_i))}{\alpha} \qquad (5)$$

According to this result, an increase in price has an impact on the sustainability effort. The sustainability effort cancels out in the case where the sustainability effort coefficient α is equal to the sustainability elasticity $K$. In this case, $s_i$ will correspond to the marginal profit.

The purchase price to be negotiated affects the sustainability effort. The store has no interest in having a purchase cost $\omega_i$, zero.

**Proof.** Maximize profit by seeking the value of $s_i$

$$\frac{d\pi_i}{ds_i} = -kc_i + kp_i - \alpha s_i - k\omega_i ;$$

Since $\frac{d^2\pi_i}{d^2s_i} = -\alpha \leq 0$ o the function admits a maximum, obtained when $\frac{d\pi_i}{ds_i} = 0$.

Since $s_i^*(p_i)$ has the equation of a line and $\frac{ds_i^*}{dp_i} = \frac{k}{\alpha} > 0$. Then the function is increasing with respect to $s_i$.

In the model shown, the store actually has two strategies to apply to improve profit. In this case, $s_i^*$, and $p_i^*$ are always profit-maximizing solutions.

**Proof.** The hessian matrix is as follows:

$$H\left[\pi_i(p_i, s_i)\right] = \begin{pmatrix} -2b & k \\ k & -\alpha \end{pmatrix}$$

$Det(H) = -k^2 + 2b\alpha$
$Tr(H) = -(2b + \alpha)$

Since $Det(H) > 0$ et $Tr(H) < 0$, then $\pi_i(p_i, s_i)$ admits a maximum in:

$$\begin{cases} p_i^* = \dfrac{1}{2}\left[\dfrac{a + ks_i}{b} + (c_i + \omega_i)\right] \\[2mm] s_i^* = \dfrac{k[p_i - (c_i + \omega_i)]}{\alpha} \end{cases}$$

Maximum profit can therefore be written as:

$$\begin{aligned} \max(\pi_i) &= -\frac{k^2 p_i^2}{2\alpha} - \frac{k^2 s_i^2}{4b} - \frac{k^3 s_i(c_i + w_i)}{2b\alpha} + \frac{1}{2b\alpha}p_i\left(k^3 s_i + k^2(a + b(c_i + w_i))\right) \\ &\quad + \frac{\alpha(a^2 + b^2 c_i^2) - 2c_i(a(k^2 + b\alpha) - b^2\alpha w_i) + w_i(-2a(k^2 + b\alpha) + b^2\alpha w_i)}{4b\alpha} \\ &\quad - (F_i + A_i y_i) \end{aligned}$$

We'll try to analyze the price that store i will negotiate with the supplier so that the store ensures maximum profit. In this case, we assume that the store is the leader.

So, we replace Eq. (5) in Eq. (4). We then obtain the negotiated price $\widehat{w_i}$:

$$\widehat{w_i} = \frac{a\alpha + (k^2 - 2b\alpha)p_i}{k^2 - b\alpha} - c_i \qquad (6)$$

*2) Supplier strategies:* In this sub-section, we focus on the supplier's profit in the supply chain. Firstly, we study the supplier's profitability in relation to store $i$, then the profitability in the whole network.

The profit function is as follows:

$$\pi_s = \sum_{i=1}^{n} \pi_{i,s} - C(\tau) + I(\nu) \qquad (7)$$

With $\sum_{i=1}^{n} d_i = D, C(\tau) = \tau.\theta D, I(\nu) = (1 - \theta)\nu D$

$C(\tau)$ corresponds to the customs fees that the supplier pays to the government for the percentage $\theta$ of the quantities imported. On the other hand, the supplier receives a government subsidy $I(\nu)$ for the percentage $(1 - \theta)$ of quantities made from local suppliers.

*a) The* supplier's strategy, considering its profit in relation to a store $i$

The supplier's profit is written as follows:

$$\pi_{i,s} = d_i(\omega_i - c - \eta_i) - \gamma\omega_i\psi_i \qquad (8)$$

Replacing Eq. (1) in Eq. (8) gives the following equation:

$\pi_{i,s} = -a(c + \eta_i) - \gamma\omega_i\psi_i + bp_i(c + \eta_i - \omega_i) + a\omega_i - ks_i(c + \eta_i - \omega_i)$

The supplier focuses solely on the pricing strategy. Always considering that the store is the leader and the supplier is the follower, we obtain the price that allows us to reach the maximum as follows:

$$\omega_i^* = \frac{(c_i + c + \eta_i)}{2} + \frac{a}{2b} + \frac{k^2(ks_i - \gamma\psi_i - 2bc_i)}{2b(k^2 + b\alpha)} \qquad (9)$$

To be able to offer a price, the supplier needs to have information on the sustainability effort of store i. It is also influenced by the cost of emissions $CO_2$ $\eta_i$.

**Proof.** To find the maximum profit, we derive the profit function $\pi_s$

$$\begin{aligned} \frac{d\pi_s}{d\omega_i} &= a + bc + \frac{ab\alpha}{k^2} + \frac{b^2 c\alpha}{k^2} - bc_i + \frac{b^2 \alpha c_i}{k^2} + ks_i + 2(-b \\ &\quad - \frac{b^2\alpha}{k^2})w_i + b\eta_i + \frac{b^2\alpha\eta_i}{k^2} \end{aligned}$$

We have: $\frac{d^2\pi_i}{d^2\omega_i} = -2(b + \frac{b^2\alpha}{k^2}) \leq 0$, then the function admits a maximum. This maximum is obtained by solving $\frac{d\pi_i}{d\omega_i} = 0$

Since $\omega_i^*(s_i)$ has the equation of a straight line and $\frac{dp_i^*}{ds_i} = \frac{k}{2b} > 0$. Then the function is increasing with respect to $s_i$.

To calculate the Supplier's profit, we first calculate the profit for each store i.

$$\frac{d\pi_i}{dp_i} = -2bp_i + a + b(c_i + \omega_i) + ks_i$$

The Eq. (9) aims to determine the wholesale price $\omega_i^*$ negotiated by the supplier in the context of sustainability, the supply chain, and operational costs. Two key parameters, $k$ (sustainability effort elasticity) and $\alpha$ (the cost coefficient of sustainability effort), play a crucial role.

Thus, $k$ measures how consumers respond to the sustainability efforts undertaken by store i. A high $k$ value indicates a strong consumer response to sustainability, meaning they are willing to purchase more sustainable products. Consequently, the supplier may consider raising the wholesale price $\omega_i^*$ without compromising demand. Consumers are willing to pay a premium for sustainable products, which can increase the supplier's profit margin. However, $\alpha$ quantifies the costs associated with the sustainability initiatives of store i. A high $\alpha$ indicates higher costs to implement sustainable practices, such as using environmentally friendly materials or reducing carbon emissions. These additional costs can exert upward pressure on the wholesale price $\omega_i^*$ negotiated by the supplier. The supplier must offset these costs to maintain profit margins.

*b) The supplier's strategy, considering the whole network:* In this section, it is necessary to focus on the overall profit of the supplier of the whole network, which will incorporate other elements not directly dependent on the store's demands. Thus, we will attempt to break down the overall profit of the supplier, as provided in Eq. (7). The total demand is expressed as follows:

$$D = an - b\sum_{i=1}^{n} p_i + k\sum_{i=1}^{n} s_i \qquad (10)$$

Based on the work of [41], it is possible to consider φ as an endogenous decision variable, with: $\varphi = \frac{p_i - w_i}{p_i}$ where $\varphi \in [0,1]$. Therefore, we can write $p_i = \frac{w_i}{(1-\varphi)}$.

Similarly, and based on the results obtained during the analysis of store-level strategies, we know that sustainability efforts $s_i$ directly impact the price offered $p_i$ by the store. Therefore, we can propose a simple linear equation to express

sustainability effort as a function of price, thus creating an equation as follows*: $s_i = \Omega p_i, \; \Omega > 0$*

When we substitute, Eq. (7) will be then expressed in the following form:

$$\pi_s = anU + (a-b)\sum_{i=1}^{n} w_i - bU \sum_{i=1}^{n} \frac{w_i}{1-\varphi} -$$
$$k\left(U\sum_{i=1}^{n} \frac{\Omega w_i}{1-\varphi} + \sum_{i=1}^{n} \frac{\Omega w_i(w_i-\eta_i)}{1-\varphi}\right) - a\sum_{i=1}^{n} \eta_i +$$
$$b\sum_{i=1}^{n} \frac{w_i\eta_i}{1-\varphi} - \gamma\sum_{i=1}^{n} w_i\psi_i \qquad (11)$$

With, $U = c - v + \theta(v - \tau)$

Based on this equation, we will attempt to examine potential strategies for the supplier, including the wholesale price and the inventory quantity reserved for store i ($\psi_i$), to improve supply chain resilience.

$$w_i{}^s = \frac{1}{2k\Omega}[(a - b - \gamma\psi_i)(1 - \varphi) + (b + k\Omega)(-U + \eta_i)]$$
$$w_i{}^s = \frac{1}{2k\Omega}[(a - b - \gamma\psi_i)(1 - \varphi) + (b + k\Omega)(-c + v - \theta(v - \tau) + \eta_i)] \qquad (12)$$

Proof. To find the maximum profit, we derive the profit function $\pi_s$

$$\frac{d\pi_s}{d\omega_i} = a - b - \frac{bU}{1-\varphi} - \frac{b\eta_i}{-1+\varphi} - k(-\frac{U\Omega}{-1+\varphi} + \frac{2\Omega w_i + \Omega\eta_i}{1-\varphi}) - \gamma\psi_i$$

We have: $\frac{d^2\pi_i}{d^2\omega_i} = -k\frac{2\Omega w_i}{1-\varphi} \le 0,$ then the function admits a maximum. This maximum is obtained by solving $\frac{d\pi_i}{d\omega_i} = 0$

*c) Resilience and wholesale price analysis:* First and foremost, it is essential to emphasize that the chosen model in this article involves the incorporation of imports to meet a portion of the store's demand from overseas. This strategic decision is made with the primary aim of enhancing the resilience of the supply chain. Consequently, it becomes crucial to examine the wholesale price concerning this resilience concept.

Now, let's explore how the wholesale price behaves within the context of government subsidies, customs fees, and the supplier's flexibility in adjusting φ to impact wholesale prices. This analysis will provide a comprehensive understanding of the dynamics at play in the supply chain and its price regulation.

The Eq. (12) represents the wholesale price offered by the supplier as a function of the level of imported quantities $w_i{}^s(\theta)$, taking into account government subsidies $v$ and custom fees $\tau$. When $v$ is high, meaning that the government offers generous subsidies for local products, the function tends to decrease with an increase in $v$ (i.e., an increase in imports). In other words, government subsidies reduce the wholesale price to encourage local production and reduce dependence on imports.

However, when $\tau$ is high, indicating substantial custom fees imposed on imports, the function $w_i{}^s(\theta)$, tends to increase with an increase in $\theta$,. Additional customs fees raise the cost of imports, which can result in higher wholesale prices for imported products.

This means that the government plays a significant role in market price regulation. Furthermore, the supplier can adjust φ to slightly lower or raise wholesale prices, but with a limited variation.

Proof. The expression for the slope of the line in the Eq. (12) is as follows: $m = \frac{1}{2k\Omega}(\tau - v)$. So, when $v > \tau$ (government subsidies exceed customs fees), $(\tau - v)$ is negative.

This results in a positive slope (m > 0), indicating a positive incline of the line.

Conversely, when $v < \tau$ (government subsidies are less than customs fees), $(\tau - v)$ becomes positive.

This yields a negative slope (m < 0), signifying a negative incline of the line

*d) Sustainability effort elasticity impact on wholesale prices:* Furthermore, an analysis regarding the elasticity of sustainability effort ($k$) highlights that with a high $k$, the function becomes more responsive to variations in $\theta$. An increase in $\theta$ (more imports) can lead to a more significant rise in wholesale prices when sustainability effort is high. This indicates that an increased commitment to sustainability can have a greater impact on the supplier's pricing decisions. On the other hand, for a low value of ($k$), indicating low elasticity of sustainability effort, the wholesale price will be less sensitive to variations in ($\theta$). An increase in ($\theta$) may have a less significant impact on wholesale prices when sustainability effort is low. In this case, other factors, such as government subsidies and customs fees, may play a more prominent role in determining wholesale prices.

*3) Government strategies:* The government derives its profit from the taxes it imposes on the quantities of products imported from abroad. In this government profit formulation, we must include not only the total quantities ordered by the stores but also the safety quantities planned by the supplier to ensure resilience. Furthermore, within the context of government profit analysis and sustainability considerations, we take into account social charges. These charges, directed towards essential programs such as healthcare and pensions, play a crucial role in bolstering stability, alleviating poverty, supporting employment, and promoting social equity, thereby contributing to a more sustainable and equitable society.

The profit equation for the government is given as follows:

$$\pi_g = [\theta(\tau + v) - v]\sum_{i=1}^{n}(d_i + \psi_i) + \lambda\frac{A}{d_m}\sum_{i=1}^{n} d_i \quad (13)$$

When examining government strategies, the focal point is the calibration of tax rates and subsidies to boost the consumption of domestic products, all the while taking into account the challenges associated with sustainability and

resilience. To conduct this analysis and streamline the study, we will rely on the wholesale price of the supplier examined in the preceding section.

Pour faciliter l'analyse, nous allons tenir compte de des subventions.

The government profit function can be written as follows:

$$\pi_g(v) = ((-1+\theta)v + \theta\tau)(an + \frac{k\Omega-b}{1-\varphi}\sum_{i=1}^{n} w_i +$$

$$\frac{n(b+k\Omega)(-c+v-\theta(v-\tau))}{2k\Omega} + \frac{(b+k\Omega)}{2k\Omega}\sum_{i=1}^{n}\eta_i + \frac{(1-\varphi)}{2k\Omega}[n(a-b) -$$

$$\gamma\sum_{i=1}^{n}\psi_i]) + \frac{aAn\lambda}{d_m} + \frac{A\lambda(-b+k\Omega)}{d_m\ (1-\varphi)}\sum_{i=1}^{n} w_i \qquad (14)$$

To simplify the analysis, we will only consider the government subsidy parameter. Indeed, addressing the government profit function has allowed us to obtain an optimal government subsidy that maximizes this profit function. So, Eq. (15) provides the government subsidy that maximizes the profit function $\pi_g(v)$.

$$v^* = \frac{n^2(c-\overline{\eta})}{2(1-\theta)} + \frac{1}{2(b+k\Omega)(1-\theta)}\left[2k\Omega\left(1 - n^2\frac{(b-k\Omega)}{-1+\varphi}\overline{w}\right) + (n^2\gamma\overline{\psi}+a-b)(1-\varphi)\right)\right] \qquad (15)$$

With,

$$\sum_{i=1}^{n}\eta_i = n\overline{\eta}, \sum_{i=1}^{n} w_i = n\overline{w}, \sum_{i=1}^{n}\psi_i = n\overline{\psi}$$

**Proof.** To find the maximum profit, we derive the profit function $\pi_g(v)$.

We have: $\frac{d^2\pi_g}{d^2 v} = -\frac{(b+k\Omega)n(-1+\theta)^2}{k\Omega} \le 0$, then the function admits a maximum. This maximum is obtained by solving $\frac{d\pi_g}{dv} = 0$

The Eq. (15) suggests that the government subsidy $v$ is determined based on a combination of factors related to the number of stores, operating costs, sustainability sensitivity, price elasticity, and various parameters associated with sustainability efforts and costs. The quotient ($\Omega$) highlights the importance of sustainability efforts relative to the price of the product, indicating that the government subsidy is influenced by the sustainability quotient.

Moreover, the government subsidy ($v$) reveals a nuanced relationship with the parameter $\theta$, representing the ratio of the quantity of foreign products in the market. The presence of ($\theta$) in the equation underscores the government's strategic approach to balancing the consumption of domestic and foreign products. As ($\theta$) increases, indicating a higher reliance on imported goods, the government subsidy adjusts to incentivize and support domestic product consumption. This reflects the government's commitment to fostering economic resilience by promoting a balance between local and international products.

## IV. NUMERICAL ANALYSIS

In this part, numerical analysis is employed to validate the conclusions drawn in the preceding section. The focus is on

exploring sustainability and resilience within logistics supply chains concerning pricing strategies and sustainability efforts. Additionally, government interventions in the logistics system are addressed, considering the constraints previously outlined in this paper. The numerical values chosen for this analysis are derived from a comprehensive examination of existing literature, ensuring alignment with established methodologies. Moreover, these values are thoughtfully selected based on the specific assumptions outlined in our study, thereby enhancing the overall validity and reliability of our numerical approach.

### A. Price, Sustainability, and Profit Analysis on the Store Side

To assess the optimal price for a store concerning sustainability parameters, the analysis highlights a positive correlation between the store's price and its sustainability effort, as illustrated in Fig. 2. This positive relationship remains consistent across different scenarios of sustainability effort elasticity ($k$), suggesting that higher values amplify the price ($p_i$) response to changes in sustainability initiatives ($s_i$). Practically, as the store strengthens its commitment to sustainability increases, there is a notable rise in product prices, with this response being particularly pronounced with higher values. Additionally, the wholesale price ($\omega_i$) plays a crucial role in determining the baseline cost and influencing the overall pricing structure.

However, in analyzing the store's profit through optimal sustainability efforts as illustrated in Fig. 3, one can also observe the correlation of sustainability ($s_i$) with the price ($p_i$) set by the store. Nevertheless, the cost coefficient of sustainability effort ($\alpha$) for store plays a crucial role in adjusting this correlation.

Furthermore, concerning the store's profit, sustainability effort elasticity ($k$) affects this profit. It can be observed in Fig. 4 that a small variation in profit is guaranteed when the sustainability effort elasticity coefficient ($k$) is small. Conversely, when this coefficient ($k$) is large, it is noticeable that it increases profit with a significant variation relative to the price.

$a = 100, b = 1, F = 30, A = 50, y = 2, \alpha = 0.2,$

$w = 40, c_i = 10, s_i = 11$



Fig. 2. Correlating Optimal store prices with sustainability parameters.

Fig. 3.   Correlating Optimal sustainability effort store price.



Fig. 4.   Influence of sustainability effort elasticity on store profit.

In the same context, the store is tasked with determining the purchase price ($p_i$) for negotiation with the supplier. Our analysis delves into the interplay between this purchase price ($p_i$) and two pivotal parameters: ($\alpha$), the cost coefficient of sustainability effort for the store, and ($k$), the elasticity of sustainability effort. This investigation is illustrated in Fig. 5, where specific values for $\alpha$ (0.01, 0.2, and 0.4) were selected to explore the nuanced dynamics of sustainability and pricing strategies.

The graph in Fig. 5 demonstrates that as sustainability sensitivity ($k$) increases, reflecting a greater commitment to sustainability, the store encounter a corresponding upward trend in optimal negotiation costs ($\widehat{w_i}$). This implies that striving for higher sustainability standards may necessitate the store to allocate additional resources to negotiation processes, potentially incurring higher expenses. The impact of this relationship varies based on the sustainability cost coefficient ($\alpha$), with lower ($\alpha$) values resulting in a relatively moderate increase in negotiation costs as sustainability sensitivity rises. This suggests that stores with lower sustainability costs may find it economically feasible to invest more in negotiations for enhanced sustainability.

*B. Partial Supplier Profit Case*

The relevant parameters are assigned as:

$$c = 10, \eta = 2, a = 50, b = 0.2, \gamma = 5, \psi = 4, c_i = 2$$

Sustainability Coefficient ($\alpha$) and Optimal Price ($\omega_i$): The graph in the Fig. 6 depicts the variation in the optimal price $\omega_i$ as a function of the sustainability coefficient ($\alpha$) for different values of the elasticity of sustainability effort ($k$). This illustrates how the optimal price of a product or service changes in response to variations in the sustainability coefficient, which can be interpreted as a measure of a company's commitment to sustainable practices. The different curves for $k$ = 0.1, 0.5 and 0.9 show how the elasticity of sustainability effort affects the sensitivity of the optimal price to changes in sustainability.

*Impact of Sustainability Effort Elasticity ($k$):* The curves reveal that the optimal price reacts differently to changes in the sustainability coefficient ($\alpha$) depending on the value of the elasticity of sustainability effort ($k$). For instance, a higher value of $k$ ($k$ = 0.9) demonstrates a greater responsiveness of the optimal price to sustainability variations compared to a lower value of $k$ ($k$ = 0.1). This suggests that, in this model, an increased commitment to sustainability (increasing $\alpha$) has a more significant impact on price when the elasticity of sustainability effort is higher.

*Profit Optimization and Sustainability:* The model appears to seek a balance between profit maximization (represented by the optimal price formula) and the promotion of sustainable practices (embodied by the sustainability coefficient $\alpha$). Price variations in response to changes in $\alpha$ and $k$ may indicate how a company can adjust its prices to achieve its economic objectives while promoting sustainability, which has significant implications for decision-making in a business context.

$$a = 100, \ b = 1, \ \Omega = 3, \ \ \gamma = 0.5, \ \phi = 0.2, \ \psi = 0.7,$$
$$c = 20, \nu = 10, \theta = 0.4, \tau = 0.8, \eta = 0.5, \alpha = 0.01$$

Exploring equilibrium wholesale prices involves an examination of the negotiated purchase price ($\widehat{w_i}$) by the store to optimize its profit and the selling price ($w_i{}^s$) set by the supplier for store $i$. The graph in Fig. 7 reveals the intersection, indicating an equilibrium price where the curves representing both prices meet. Additionally, we underscore the significance of the sustainability sensitivity coefficient ($k$) in this analysis.



Fig. 5.   Comparative analysis of optimal negotiation price and sustainability parameters.

Fig. 6.  Wholesalse price in function of sustainability coefficient $k$.



Fig. 7.  Equilibrium wholesale prices.

## C. Supplier Chain Resilience

In an effort to investigate the resilience of the supply chain, we have constructed a set of numerical values given in Table II. The objective is to analyze the profits of both the government and the supplier. The numerical data within the table offers a comprehensive view of how variations in parameters, such as the rate of foreign products ($\theta$) and inventory security sensitivity ($\gamma$), impact the financial outcomes for both key stakeholders in the supply chain. This analysis provides valuable insights into the resilience of the supply chain under different conditions and aids in decision-making related to supply chain management strategies.

TABLE II.  IMPACT OF SUSTAINABILITY AND RESILIENCE PARAMETERS ON GOVERNMENT AND SUPPLIER PROFITS

| $\theta$ | $\lambda$ | $\pi_s$ | $\pi_g$ |
|---|---|---|---|
| 0.39 | 0.3 | 1101.00 | 40291 |
| 0.39 | 0.54 | 1065.29 | 38371 |
| 0.39 | 0.78 | 1029.59 | 36451 |
| 0.59 | 0.3 | 5530.22 | 39956.7 |
| 0.59 | 0.54 | 5368.52 | 38036.7 |
| 0.59 | 0.78 | 5206.82 | 36116.7 |
| 0.79 | 0.3 | 9981.19 | 39622.4 |
| 0.79 | 0.54 | 9693.49 | 37702.4 |

## V.    RESULTS AND DISCUSSION

In this section, our attention is directed towards a comprehensive examination of the outcomes derived from the implemented model and numerical analyses. Subsequently, the ensuing discourse will meticulously explore the findings pertaining to profit maximization for each stakeholder, specifically, the store, the supplier, and the government.

Within this framework, stores wield authority over pricing and sustainability strategies, while the supplier assumes responsibility for distribution and inventory management. Simultaneously, the government plays a crucial role by enforcing regulations, particularly through taxation and subsidies, with an emphasis on both domestic and foreign products. Central to our investigation is the resilience of this supply chain, where we take into account various factors such as security stock, the social impact of stores, and the diversification of supply sources. Our primary objective was to comprehensively analyze the dynamic interplay within this supply chain, with a particular emphasis on sustainability, resilience and the promotion of domestic production.

Each stakeholder of the supply chain; namely the store, the supplier, and the government, adopt strategies to maximize their profits. Regarding the store i, the model presents two strategies for improving profit, with $s_i^*$ and $p_i^*$ as always profit-maximizing solutions. From the supplier's perspective, the first strategy centers on calibrating the optimal wholesale price ($\omega_i^*$) in a manner that harmonizes sustainability commitments with supply chain. The second strategy expands this focus to encompass the broader network, seeking to establish a wholesale price ($w_i^s$) that incorporates resilience metrics, thus ensuring profit maximization across the network. On the governmental front, the strategies revolve around custom fees ($\tau$) on international imports and subsidies ($\nu$) to bolster the competitiveness of domestic products, balancing international trade with local economic encouragement.

Sustainability emerges as a key point in the main conclusions of this study, particularly in its influence on pricing strategies within supply chains. Firstly, a store's pricing strategy is closely linked to its sustainability efforts ($s_i$), with the cost coefficient of these efforts ($\alpha$) playing a significant role in shaping pricing ($p_i$). The elasticity of sustainability efforts ($k$) is found to directly impact profit margins ($\pi_i$) where a lower elasticity results in smaller profit variations, and a higher elasticity leads to more substantial profit fluctuations relative to price changes. Additionally, as a store intensifies its commitment to sustainability, it incurs higher negotiation costs ($\widehat{w_i}$), though this is less burdensome for stores with lower associated sustainability costs ($\alpha$). The optimal price ($\omega_i^*$) of a product is thus affected by the company's sustainability commitment ($\alpha$) and its responsiveness to sustainability efforts ($k$), with more responsive companies experiencing more significant pricing effects. This indicates a strategic imperative for businesses to align pricing with sustainability objectives, balancing environmental considerations with the aim of profit optimization.

Resilience in the model is a keystone of the supplier's strategy, by allocating reserved inventory quantities ($\psi_i$) in

wholesale prices for each store i. This approach ensures that pricing and stock level decisions are in concert with the overarching ambition to fortify the supply chain's robustness. Parallel to this, the government's regulatory role is crucial in preserving a delicate balance between encouraging domestic products and regulating the import of foreign products, thereby providing more flexibility to diversify the product sources and underpinning economic resilience. Indeed, the resilience of the supply chain, as examined through numerical analysis, is affected by parameters like the rate of foreign products ($\theta$) and inventory security sensitivity ($\gamma$), which in turn influence the profits of both the government and the supplier.

The equilibrium wholesale price is found at the intersection of the store's purchase price and the supplier's selling price, with the sustainability sensitivity coefficient being a significant factor in this determination. This convergence underscores the indispensable importance of collaboration among supply chain stakeholders in achieving optimal profits. Through a synergistic partnership, suppliers, stores, and government entities can refine sustainability efforts.

## VI. Conclusion

Sustainability and resilience have become essential pillars in the formulation of pricing policies, ensuring that economic strategies are adaptive and viable over the long term. These concepts, grounded in economic significance as demonstrated by the literature, require in-depth analysis of how pricing policies and governmental regulatory measures, such as the promotion of domestic products, can coexist to strengthen the sustainability and resilience of supply chains.

This study addresses the under-researched interplay between sustainability, resilience, and the promotion of domestic products within supply chains, offering new approaches for tackling the complex challenges in supply chain management. It highlights the unexpected competitive dynamics that can occur even for monopolistic suppliers when they supply and compete with their stores. Our supply chain management model scrutinizes the intricate interactions among three pivotal entities: a central supplier, multiple stores, and the government.

The supply chain framework is analyzed with a focus on the roles of stores, suppliers, and the government. Stores control pricing and sustainability strategies, suppliers manage distribution and inventory, and the government enforces regulations, including taxation and subsidies. The study emphasizes the resilience of the supply chain, considering factors like security stock, social impact, and diversification of supply sources. Stakeholders adopt profit-maximizing strategies, with stores having two pricing strategies linked to sustainability efforts. Suppliers focus on optimal wholesale prices aligned with sustainability and resilience metrics. Government strategies involve custom fees and subsidies to balance international trade and support domestic products. Sustainability efforts significantly influence pricing strategies, with lower elasticity resulting in smaller profit variations. Resilience is crucial for suppliers, involving reserved inventory quantities. Government regulation balances encouraging domestic products with regulating imports, enhancing economic resilience. The equilibrium wholesale price is

determined by collaboration among stakeholders, emphasizing the importance of a synergistic partnership for optimal profits and sustainability efforts.

It is crucial to elucidate the limitations of our study, providing researchers with valuable context. The main limitation of the model developed in this research lies in its assumption that demand is deterministic, a simplification that does not reflect the dynamic and often unpredictable reality of the market. In practice, consumer demand is subject to fluctuations influenced by various factors such as changing preferences, competition, and economic conditions. This deterministic approach can lead to inaccurate forecasts and suboptimal decisions in a real business context. To enhance the practical relevance of this study, it would be crucial to explore more sophisticated models that incorporate demand variability, allowing companies to better adapt to market changes and optimize their strategies by considering inherent uncertainties.

Looking ahead, a prospective avenue for upcoming research studies involves extending the current model to incorporate a more complex network structure. This expansion would enable a more in-depth exploration of the cooperative dynamics among supply chain stakeholders. A particular focus should be placed on developing a stochastic model to better align with the real-world scenario where demand is variable, allowing for a more accurate representation of uncertainties. The objective is to examine whether the observed impacts in this study persist or if new patterns emerge within the intricacies of a larger supply chain network. Such an investigation could significantly enrich our understanding of supply chain sustainability, resilience, and collaboration by providing insights that are more reflective of the complexities and uncertainties inherent in real-world demand dynamics.

## Conflicts of Interest

Authors declare no conflict of interest.

## References

[1] D. Saidi, J. E. Alami, et M. Hlyal, « Building Sustainable Resilient Supply Chains in Emerging Economies: Review of Motivations and Holistic Framework », IOP Conf. Ser.: Earth Environ. Sci., vol. 690, no 1, p. 012057, mars 2021, doi: 10.1088/1755-1315/690/1/012057.

[2] D. Saidi, J. El Alami, et M. Hlyal, « Sustainable Supply Chain Management: review of triggers, challenges and conceptual framework. », in IOP Conference Series: Materials Science and Engineering, IOP Publishing, 2020, p. 012054.

[3] D. Saidi, K. Jharni, J. Alami, et M. Hlyal, « WHAT MODELING APPROACHES USED FOR A SUSTAINABLE RESILIENT SUPPLY CHAIN », Journal of Theoretical and Applied Information Technology, vol. Vol.100, p. 35, déc. 2022, doi: 10.5281/zenodo.7542949.

[4] A. Ranjan et J. K. Jha, « Pricing and coordination strategies of a dual-channel supply chain considering green quality and sales effort », Journal of Cleaner Production, vol. 218, p. 409-424, 2019, doi: 10.1016/j.jclepro.2019.01.297.

[5] H. Rajabzadeh et R. Babazadeh, « A game-theoretic approach for power pricing in a resilient supply chain considering a dual channel biorefining structure and the hybrid power plant », Renewable Energy, vol. 198, p. 1082-1094, 2022, doi: 10.1016/j.renene.2022.08.118.

[6] A. De et S. P. Singh, « A resilient pricing and service quality level decision for fresh agri-product supply chain in post-COVID-19 era », International Journal of Logistics Management, vol. 34, no 4, p. 1101-1140, 2023, doi: 10.1108/IJLM-02-2021-0117.

[7] S. Amiri-Pebdani, M. Alinaghian, et S. Safarzadeh, « Time-Of-Use pricing in an energy sustainable supply chain with government interventions: A game theory approach », Energy, vol. 255, p. 124380, sept. 2022, doi: 10.1016/j.energy.2022.124380.

[8] A. Jafari-Nodoushan, M. H. D. Sadrabadi, M. Nili, A. Makui, et R. Ghousi, « Designing a sustainable disruption-oriented supply chain under joint pricing and resiliency considerations: a case study », Computers & Chemical Engineering, p. 108481, oct. 2023, doi: 10.1016/j.compchemeng.2023.108481.

[9] I. E. Nielsen, S. Majumder, S. S. Sana, et S. Saha, « Comparative analysis of government incentives and game structures on single and two-period green supply chain », Journal of Cleaner Production, vol. 235, p. 1371-1398, oct. 2019, doi: 10.1016/j.jclepro.2019.06.168.

[10] J.-Y. Chen, S. Dimitrov, et H. Pun, « The impact of government subsidy on supply Chains' sustainability innovation », Omega, vol. 86, p. 42-58, juill. 2019, doi: 10.1016/j.omega.2018.06.012.

[11] A. Toptal, H. Özlü, et D. Konur, « Joint decisions on inventory replenishment and emission reduction investment under different emission regulations », International Journal of Production Research, vol. 52, no 1, p. 243-269, janv. 2014, doi: 10.1080/00207543.2013.836615.

[12] B. Zhang et L. Xu, « Multi-item production planning with carbon cap and trade mechanism », International Journal of Production Economics, vol. 144, no 1, p. 118-127, juill. 2013, doi: 10.1016/j.ijpe.2013.01.024.

[13] X. Chen, X. Wang, V. Kumar, et N. Kumar, « Low carbon warehouse management under cap-and-trade policy », Journal of Cleaner Production, vol. 139, p. 894-904, déc. 2016, doi: 10.1016/j.jclepro.2016.08.089.

[14] D. Krass, T. Nedorezov, et A. Ovchinnikov, « Environmental Taxes and the Choice of Green Technology », Production and Operations Management, vol. 22, no 5, p. 1035-1055, 2013, doi: 10.1111/poms.12023.

[15] B. Yalabik et R. J. Fairchild, « Customer, regulatory, and competitive pressure as drivers of environmental innovation », International Journal of Production Economics, vol. 131, no 2, p. 519-527, juin 2011, doi: 10.1016/j.ijpe.2011.01.020.

[16] Z. Yanju, H. Fengying, et Z. Zhenglong, « Study on joint contract coordination to promote green product demand under the retailer-dominance », Journal of Industrial Engineering and Engineering Management, vol. 34, no 2, p. 194-204, 2020, doi: 10.13587/j.cnki.jieem.2020.02.021.

[17] K. Yan, G. Hua, T. Cheng, T. Choi, J. Dong, et X. Li, « Optimal Pricing and Quality Decisions Under Cooperative Promotion of Cross-Market Service Platforms », IEEE Transactions on Engineering Management, p. 1-25, 2023, doi: 10.1109/TEM.2023.3301886.

[18] H. Yang et W. Chen, « Retailer-driven carbon emission abatement with consumer environmental awareness and carbon tax: Revenue-sharing versus Cost-sharing », Omega, vol. 78, p. 179-191, juill. 2018, doi: 10.1016/j.omega.2017.06.012.

[19] Q. Li et B. Shen, « Sustainable Design Operations in the Supply Chain: Non-Profit Manufacturer vs. For-Profit Manufacturer », Sustainability, vol. 8, no 7, Art. no 7, juill. 2016, doi: 10.3390/su8070639.

[20] Z. H. Adnan, K. Chakraborty, S. Bag, et J. S. Wu, « Pricing and green investment strategies for electric vehicle supply chain in a competitive market under different channel leadership », Annals of Operations Research, 2023, doi: 10.1007/s10479-023-05523-y.

[21] A. M. Ledari, A. A. Khamseh, et B. Naderi, « Pricing Models for a Two-echelon Supply Chain with Substitute and Complementary Products Considering Disruption Risk », International Journal of Supply and Operations Management, vol. 10, no 2, p. 187-208, 2023, doi: 10.22034/ijsom.2022.108726.1903.

[22] A. A. Taleizadeh, A. Ghavamifar, et A. Khosrojerdi, « Resilient network design of two supply chains under price competition: game theoretic and decomposition algorithm approach », Oper Res Int J, vol. 22, no 1, p. 825-857, mars 2022, doi: 10.1007/s12351-020-00565-7.

[23] G. Lyu, M. Zhao, Q. Ji, et X. Lin, « Improving resilience via capacity allocation and strategic pricing: Co-opetition in a shipping supply chain », Ocean and Coastal Management, vol. 244, 2023, doi: 10.1016/j.ocecoaman.2023.106779.

[24] S. P. Parvasi, A. A. Taleizadeh, et L. E. Cárdenas-Barrón, « Retail price competition of domestic and international companies: A bi-level game theoretical optimization approach », RAIRO - Operations Research, vol. 57, no 1, p. 291-323, 2023, doi: 10.1051/ro/2023007.

[25] S. P. Parvasi et A. A. Taleizadeh, « Competition pricing between domestic and foreign manufacturers: a bi-level model using a novel hybrid method », Sadhana - Academy Proceedings in Engineering Sciences, vol. 46, no 2, 2021, doi: 10.1007/s12046-021-01627-y.

[26] A. Mostofi, V. Jain, Y. Mei, et L. Benyoucef, « A new pricing mechanism for pharmaceutical supply chains: a game theory analytical approach for healthcare service », International Journal of Logistics Research and Applications, 2022, doi: 10.1080/13675567.2022.2122421.

[27] H. Rajabzadeh, A. Arshadi Khamseh, et M. Ameli, « A Game-Theoretic Approach for Pricing Considering Sourcing, Andrecycling Decisions in a Closed-Loop Supply Chain Under Disruption », Communications in Computer and Information Science, vol. 1458 CCIS, p. 137-157, 2021, doi: 10.1007/978-3-030-89743-7_9.

[28] K. Wang, Y. Zhao, Y. Cheng, et T.-M. Choi, « Cooperation or Competition? Channel Choice for a Remanufacturing Fashion Supply Chain with Government Subsidy », Sustainability, vol. 6, no 10, Art. no 10, oct. 2014, doi: 10.3390/su6107292.

[29] N. Droste et al., « Steering innovations towards a green economy: Understanding government intervention », Journal of Cleaner Production, vol. 135, p. 426-434, nov. 2016, doi: 10.1016/j.jclepro.2016.06.123.

[30] S. S. Sana, « Price competition between green and non green products under corporate social responsible firm », Journal of Retailing and Consumer Services, vol. 55, p. 102118, juill. 2020, doi: 10.1016/j.jretconser.2020.102118.

[31] D. Yang et T. Xiao, « Pricing and green level decisions of a green supply chain with governmental interventions under fuzzy uncertainties », Journal of Cleaner Production, vol. 149, p. 1174-1187, avr. 2017, doi: 10.1016/j.jclepro.2017.02.138.

[32] P. He, Y. He, et H. Xu, « Channel structure and pricing in a dual-channel closed-loop supply chain with government subsidy », International Journal of Production Economics, vol. 213, p. 108-123, juill. 2019, doi: 10.1016/j.ijpe.2019.03.013.

[33] S. Amiri-Pebdani, M. Alinaghian, et H. Khosroshahi, « Pricing in competitive energy supply chains considering government interventions to support CCS under cap-and-trade regulations: A game-theoretic approach », Energy Policy, vol. 179, 2023, doi: 10.1016/j.enpol.2023.113630.

[34] M. Rasti-Barzoki et I. Moon, « A game theoretic approach for analyzing electric and gasoline-based vehicles' competition in a supply chain under government sustainable strategies: A case study of South Korea », Renewable and Sustainable Energy Reviews, vol. 146, p. 111139, août 2021, doi: 10.1016/j.rser.2021.111139.

[35] D.-H. Lee et J.-C. Yoon, « Decisions on Pricing, Sustainability Effort, and Carbon Cap under Wholesale Price and Cost-Sharing Contracts », Sustainability (Switzerland), vol. 14, no 8, 2022, doi: 10.3390/su14084863.

[36] R. W. Anderson, « Some theory of inverse demand for applied demand analysis », European Economic Review, vol. 14, no 3, p. 281-290, janv. 1980, doi: 10.1016/S0014-2921(80)80001-8.

[37] B. Niu, J. Li, J. Zhang, H. K. Cheng, et Y. (Ricky) Tan, « Strategic Analysis of Dual Sourcing and Dual Channel with an Unreliable Alternative Supplier », Production and Operations Management, vol. 28, no 3, p. 570-587, 2019, doi: 10.1111/poms.12938.

[38] S. Y. Tang et P. Kouvelis, « Supplier Diversification Strategies in the Presence of Yield Uncertainty and Buyer Competition », M&SOM, vol. 13, no 4, p. 439-451, oct. 2011, doi: 10.1287/msom.1110.0337.

[39] S. R. Madani et M. Rasti-Barzoki, « Sustainable supply chain management with pricing, greening and governmental tariffs determining strategies: A game-theoretic approach », Computers and Industrial Engineering, vol. 105, p. 287-298, 2017, doi: 10.1016/j.cie.2017.01.017.

[40] A. Ranjan et J. K. Jha, « Pricing and coordination strategies of a dual-channel supply chain considering green quality and sales effort », Journal of Cleaner Production, vol. 218, p. 409-424, mai 2019, doi: 10.1016/j.jclepro.2019.01.297.

[41] C. Wang, W. Wang, et R. Huang, « Supply chain enterprise operations and government carbon tax decisions considering carbon emissions », Journal of Cleaner Production, vol. 152, p. 271-280, mai 2017, doi: 10.1016/j.jclepro.2017.03.051.

# Investigating Agile Values and Principles in Real Practices

Abdullah A H Alzahrani

Department of Computers, Engineering and Computing College at Alqunfuda,
Umm Al Qura University, Makkah, Saudi Arabia

*Abstract*—Software engineering is the field of development of information systems. However, the development process can often be complicated. Therefore, many researchers have introduced their approaches to manage the complication. This led to the introduction of new subfields such as change management, and organisational change. Agile can be regarded as a collection of best practices with the same values and principles. Since the introduction of Agile manifesto, many researchers, manufacturers, and organisations have introduced their thoughts, tools, and models to enhance the understanding and adoption of Agile. Sharing a similar understanding of Agile among people involved is essential in order to adopt it. This paper investigates the understanding of Agile among IT professionals. In addition, the factors that impact the understanding and adoption of Agile are highlighted and studied. A survey methodology was employed in this research among IT professionals from different organisations. The results of this study show that productivity and ability to accept change are conflicting the understanding among participants. Furthermore, the experience of participants has an impact on the ways in which Agile are adopted.

*Keywords—Agile; software engineering; information systems; change management; organisational change*

## I. INTRODUCTION

Since 2001, Agile [1]–[3] has been known to be the most adaptive way in the field of software engineering. It refers to a software development model that should accept changes during software development. Hesselberg [4] articulated that "The agile mindset is now finding its way into the C-suite, and it is starting to radically change the way organizations are led and managed. Business agility is on everybody´s lips, for very good reasons".

Some researchers define Agile as a collection of best practices with same values and principles [5]. However, other researchers define Agile as a subset of iterative methods of the traditional methodologies of software development [6]. In general, Agile can be defined as development way that relies on the philosophy of change embracing.

Agile has several benefits that overcome traditional ways. These benefits can be summarized in the following: change embracing, customer heard, quick achievement, good interactions, and continuous improvements. However, several drawbacks come along with Agile. These drawbacks can be summarized in the following: time consuming, unsatisfactory documentation, change dilemma, and unclear customer [7]–[9].

Many have introduced their Agile methods and frameworks [10]–[13] such as Scrum [14]–[16], eXtreme Programming (XP) [17]–[19], and DevOps [20]–[22]. Each method and framework has its own pros and cons. However, there is no unified framework or method that can be considered to be the best practice in every circumstance [23].

Since Agile introduction, many software development teams are claiming that they adopt Agile model, However, a question of "do they adopt agile?" can be raised. This research aims to investigate the understanding of Agile among software development teams. It is crucial that development teams share the same understanding of Agile when it is employed. Therefore, this research questions the understanding by investigation each value and principle of Agile against the participants views of these values and principles. This is carried out by asking the participants to priorities and criticizes the values and principles. In addition, it highlights any modified version of adopting and understanding that can exist.

The remainder of this paper is structured as follows: Section II describes the related work of this field of research. This is followed by the Section III which describes the research questions of this research. The methodology employed is described in Section IV. Section V is divided into three subsections which show and discuss the findings of this research. Finally, the conclusion and limitations will be drawn in Section VI.

## II. RELATED WORK

Ozkan et al. [24] have introduced a study that combines Agile principles from different resources and divided them into groups in order to develop a better understanding of Agile. This grouping was done by one expert. The evaluation process involved two experts. However, authors still see the understanding of Agile as a challenging process [25].

Nurdiani et al. [26] introduced their methodology to understand and compare Agile Maturity Models (AMM) and its strategies. The methodology was based on collecting data from previous studies' results on the topic and a survey done on 46 participants. However, the attempt arrived to transfer Agile methods onto organizational level and practical implications [27].

Koi-Akrofi et al. [28] investigated Agile in management of IT projects. The study was to compare the use of Agile and traditional ways to manage IT projects. The study focused on the challenges of using Agile. The authors found that despite the benefits of using Agile, many challenges accompany the

employment of Agile. These challenges, such as empowerment and organization culture, make it difficult to apply. Therefore, many go to the traditional ways instead of Agile or merge both models in real life.

Hess et al. [29] investigate the ways to improve the understanding of Agile in the perspective of Information needs and communication and collaboration. The study is based on a previous study of the authors that investigate traditional ways. The authors are comparing the results from both studies. The findings highlighted the gap in product inconsistency when employing Agile. However, authors found also that teams do not share the same understanding of the values of Agile as they use them differently [30].

Barroca et al. [31] introduced a paper that summarizes an international workshop discussion on Agile transformation. Many definitions were discussed and presented. In addition, challenges of Agile transformation were identified. One of the challenges is the understanding of Agile within the teams. However, the authors attempted combine solutions to overcome these challenges [32]–[34].

Jia et al. [35] have conducted a case study to investigate the understanding of developers of Agile requirements. The study was conducted on around 130 students divided into 17 teams to develop a web-based email management system. The main findings of the study identified difficulties in understanding the Agile requirements. Inconsistent understanding of Agile was noticed [36].

Baham et al. [37] have introduced a theoretical core that they found to be a gap in the studies of Agile. The authors offered a framework that unified the theoretical understanding of Agile. However, the study is considered to be an inspiration for future discussion and implications on the topic.

Eilers et al. [38] have investigated the gap between being Agile and going Agile among development teams. The study included around 129 participants and shows that the empowerment of the development teams enhances the Agile work and overcomes challenges. In addition, happiness and commitment are the factors of connection of being or doing Agile [39]. However, it has been argued that empowerment has limited impact on project outcomes [40].

## III. RESEARCH QUESTION

This research outlines two research questions which focus on the investigation of Agile understanding and adoption among software development teams.

RQ1. Is there a difference in understanding and adoption between experienced IT professionals and less experienced professionals? In order to answer this question, participants will be asked to categorize themselves into one of five categories. Based on the answer to this question, the collected data will be analyzed accordingly.

RQ2. Based on the experience of IT professionals, is there an evolutionary understanding of Agile values and principles? The answer to this question will rely on participants' views on priority and criticism of Agile values and principles.

## IV. METHODOLOGY

This research employed the methodology of questionnaire. This methodology will allow collected data based on experience of participants and classify the responses for further investigation. The questionnaire was sent to potential participants in the IT field via emails and social communication. Fig. 1 shows the stages of the investigation.

The questionnaire was sent to 80 potential participants, the responses received were 38 responses. The participants are IT professionals with different job titles namely: Software Analyst, Software Designer, Software Developer/ Implementer/ Programmer, Software Engineer, Software Project Coordinator, Software Project Manager, and Software Tester. It is clear from the aforementioned job titles that the focus on this study is Agile adoption and understanding among software development teams.



Fig. 1. Research methodology.

The questionnaire was divided into four sections. The first section includes direct questions about gender, job title, qualification, and years of experience. In addition, in this section, participants are asked to rate their knowledge on Agile and whether they have been taught or trained on Agile. Furthermore, participants are asked if they view their organization employing Agile. The second section is related to Agile values shown in Table I. The participants are asked to prioritize the values from their experience.

TABLE I. AGILE VALUES [3]

| Values | abbreviation |
|---|---|
| Individuals and interactions over processes and tools | Val1 |
| Working software over comprehensive documentation | Val2 |
| Customer collaboration over contract negotiation | Val3 |
| Responding to change over following a plan | Val4 |

In the third section of the questionnaire, the participants are asked to criticize the Agile values shown in Table I with one of the options shown Table II. The fourth section of the questionnaire is also regarding participants' criticism on the Agile principles shown in Table III. The criticism is based on participants' selection of options from Table II to the principles shown in Table III.

TABLE II. CRITICISM CRITERIA

| Criterion | Meaning |
|---|---|
| Keep as it is | No change suggested |
| Need To be Removed | Removal suggested |
| Need To Be Modified | Some changes might be suggested |

Following the ethical manner of scientific research, participants were informed that the collected data is confidential and is used for research purposes. Thereafter, they were asked to provide their consent to participate, and they were able to withdraw at any stage of the process. The participants were assured that their privacy is protected and respected.

TABLE III. AGILE PRINCIPLES [3]

| Principles | abbreviation |
|---|---|
| Our highest priority is to satisfy the customer through early and continuous delivery of valuable software. | PPl1 |
| Welcome changing requirements, even late in development. Agile processes harness change for the customer's competitive advantage. | PPl2 |
| Deliver working software frequently, from a couple of weeks to a couple of months, with a preference to the shorter timescale. | PPl3 |
| Business people and developers must work together daily throughout the project. | PPl4 |
| Build projects around motivated individuals. Give them the environment and support they need, and trust them to get the job done. | PPl5 |
| The most efficient and effective method of conveying information to and within a development team is face-to-face conversation. | PPl6 |
| Working software is the primary measure of progress. | PPl7 |
| Agile processes promote sustainable development. The sponsors, developers, and users should be able to maintain a constant pace indefinitely. | PPl8 |
| Continuous attention to technical excellence and good design enhances agility. | PPl9 |
| Simplicity--the art of maximizing the amount of work not done--is essential. | PPl10 |
| The best architectures, requirements, and designs emerge from self-organizing teams. | PPl11 |
| At regular intervals, the team reflects on how to become more effective, then tunes and adjusts its behavior accordingly. | PPl12 |

From Fig. 2, 74% of participants are males where the remaining are females. In addition, participants are from a different range of experiences. However, most of them have less than 10 years of experience with a percentage around 61%. Participants with more than 10 years of experience constitute around 39% of all participants.



Fig. 2. Participants genders and experience range.

## V. RESULTS AND DISCUSSION

In this section, the results of the research will be shown and discussed. First, the general findings of the research will be presented and discussed. These are related to the participants understanding of Agile and the source of knowledge learnt

Agile from. In addition, how the participants view the adoption of Agile in their organisations. Next, the remaining subsections are presenting the findings related to the research questions of the experiences of the participants and their understanding of Agile and new models of understanding.

### A. General findings

In order to identify the source of the participant's knowledge of Agile, participants were asked direct questions specifying if they have learned Agile in school or at training in a workplace. Fig. 3 shows that around 74% of participants have been trained in the workplace on Agile. On the other hand, just above half of participants have been taught Agile at school.



Fig. 3. Participants knowledge source of agile.

From Fig. 3, it is obvious that workplaces are more interested in Agile, and they tend to train professionals on Agile even if they have studied it. From this, organisations of participants urge the adoption of Agile. However, contrastingly, Fig. 4 illustrates the view of the participants on the adoption of Agile, and it is obvious that around 58% of participants believe that their organisations adopt Agile.



Fig. 4. Participants views on adoption of Agile in their organisation.

Fig. 5 illustrates a key point about the understanding of Agile among IT professionals. It shows the participants' self-rating on Agile knowledge. Although, as shown in Fig. 3, 74% of participants received training on Agile and 55% of them studied it in school, Fig. 5 shows that around 61% of participants rated their knowledge in Agile as medium as or less than medium.



Fig. 5.   Self-rating of participants on Agile knowledge.

It can be concluded from the general findings that there is a gap in sharing the understanding of Agile among IT professionals. In addition, despite the urge of the organisations to adopt Agile, a struggle can be noticed to do such. This might be attributed to the gap of Agile understanding among the IT professionals.

*B.  Prioritising Agile Values*

In order to investigate the new understanding of Agile, participants were asked to prioritize the values of Agile shown in Table I. It is obvious from Fig. 6 that Val2 has been prioritized the most by participants as it occupies the first and second priority. Val2 is regarding the value of working software over documentation. This is an interesting point as this value needs interactions with clients and might lead to a change of requirements. Therefore, Val1 and Val4 are following in the priority order as shown in Fig. 6.



Fig. 6.   Participants prioritising Agile values.

It can be concluded that the understanding of the participants of Agile values has been impacted by the tendencies of productivity and customer satisfaction. In addition, software documentation might be affected by adopting Agile as this understanding imposes more productivity over quality and more changes over quality.

*C.  Criticism of Agile Values and principles*

Participants were asked to provide an abstract view on the criticism of the values and principles of Agile. Fig. 7 shows the participants responses to a criticism question on each value of Agile. The question aims to collect a general answer of participants as if they believe that a value should be eliminated, modified, or kept as is.

From Fig. 7, it is obvious that Val2 is a controversial value among participants. In addition, a high number of opinions regarding modification and elimination to the value were focused on Val2 with around 71%. It is worth noting that Val2 is about productivity over quality. Furthermore, Val4 received an equal number of responses to modification and elimination, however, it received the highest number of responses to be kept as is among other Agile values.



Fig. 7.   Participants criticising Agile values.

With regards to the criticism of Agile principles, Fig. 8 illustrates the responses of the participants. Overall results show that participants tend to have no criticism of Agile principles, as the dominant response is to keep as is. However, PPL2, PPL6, PPL7, and PPL10 seem to receive responses regarding modification of these principles. It is worth noting that PPL2 is about changes in software, PPL6 is about communication with clients, PPL7 is about productivity, and PPL10 is about simplicity of software.

It can be concluded from results shown in Fig. 7 and Fig. 8 that Val2, which is regarding productivity over quality, is again obvious in the criticism of the participants. Interestingly, the principles of PPL2, PPL7, and PPL10 are related to the Val2. From this it can be noticed that productivity over software quality is a controversial understanding among the participants. In addition, this is related directly to the adoption of Agile as it might introduce new ways of Agile adoption or affect the traditional way of Agile adoption.

Fig. 8. Participants criticising Agile principles.

## D. *Criticism of Agile Values and principles Based on experience.*

In this section, the previous criticism of Agile values and principles is investigated further. The focus is on the comparison of participants' responses to the criticisms based on years of experience. The responses were divided into two groups: less than 10 years of experience and more than 10 years of experience.

Fig. 9 illustrates the criticism responses of participants regarding Agile values based on experience of participants. It is obvious from the Figure, that participants with less than 10 years of experience have no tendency to criticize except for Val2 where around 70% of the participants think Val2 needs to be eliminated or modified.



Fig. 9. Based on experience participants criticising Agile values.

On the other hand, participants with more than 10 years of experience agree with participants with less than 10 years of experience on that Val2 needs to be eliminated or modified. Furthermore, the majority of participants with more than 10 years of experience think that Val3 and Val4 should be modified and eliminated respectively.

With regards to Agile principles, Fig. 10 shows the criticism responses of participants based on experience of participants. The majority of participants with over 10 years of experience think that PPL7 should be eliminated. In addition, they think that PPL6 should be eliminated or modified. In addition, PPL10 received great attention to be modified by the majority of participants with over 10 years of experience.

On the other hand, the majority of participants with less than 10 years of experience believe that PPL2 should be eliminated or modified. In addition, they think PPL6 should be modified. Furthermore, participants with less than 10 years of experience give the same attention as participants with over 10 years of experience on that PPL10 should be modified.



Fig. 10. Based on experience participants criticising Agile principles.

It can be concluded from the results that experience impacts the understanding and perspective of participants on Agile values and principles. This is clear from the difference of the views of Val3 and Val4 as participants with over 10 years of experience are more likely to criticize these two values of Agile unlike participants with less than 10 years of experience. However, all participants agree on the criticism of Val2. In summary these different views might affect the adoption of Agile within teams in organisations.

## VI. CONCLUSION AND LIMITATIONS

In this paper, Agile values and principles were focused on in the perspective of the understanding and the adoption from

the point of view of the IT professional involved. The methodology of survey was employed in order to investigate the understanding and adoption of Agile among IT professionals.

Findings of this research can be summarized as follows:

*1) In* general, there is a clear difference in understanding and adoption of Agile.

*2) Workplaces* are taking Agile seriously and as they pay great attention to Agile training with 74% of participants been trained in their workplace. It seems that workplaces do not rely on the members knowledge of Agile from members' study in schools.

*3) Despite* the workplace attention to Agile adoption, a great deal of participants does not believe that their workplace is adopting Agile.

*4) Despite* the training they received and the school teaching of Agile, participants are not confident about their knowledge and understanding of Agile as 61% of them seem to rate their knowledge and understanding as medium of less.

*5) The* participants understand Agile as productivity more than change embracement. This is due to the priority they give to Val2 as first and second priority.

*6) However,* Val2 is a controversial value of Agile as 71% of responses in the criticism of Agile values goes to need modification or elimination of Val2.

*7) Agile* principles PPL 2, 6, 7, and 10 are criticized by participants to be modified or eliminated, however, participants in general do not criticize other Agile principles.

*8) Participants* with less than 10 years of experience tend to avoid criticizing Agile values except Val2 which they believe should be eliminated or modified.

*9) Participants* with over 10 years of experience tend to criticize Agile values in particular Val 2, 3, and 4.

*10)A* difference in focus can be seen in Agile principles criticism between participants with over 10 years of experience and participants with less than 10 years of experience. PPL2 is the principle that participants with less than 10 years of experience believe should be eliminated or modified, whereas participants with over 10 years of experience believe that PPL7 should be eliminated or modified.

As limitations of this study, the number of participants responses is one of the limitations. Reaching a higher number might help in generalizing the results and might give other perspectives to the issues of Agile understanding. Another limitation related to the number of participants is that as the number is not large, responses cannot be divided based on teams with the same job titles.

As future work based on this research, the investigation of refinement of the values and principles of agile should be conducted, since this research findings motivate the refinement of them. Another future direction is the introduction of a new model for training people in organizations which tend to adopt Agile. Finally, further investigation needs to be conducted in

the field of quality of software which is developed with the employment Agile.

## REFERENCES

[1] S. Hayes and M. Andrews, "An introduction to agile methods," Steve Hayes Khatovar Technol. Steve Khatovartech Com Httpwww Khatovartech Com, 2003, Accessed: Oct. 31, 2023. [Online]. Available: https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=0e0b d8299ccad16526f18e7a9003b8b49d1b273e

[2] K. Beck et al., "Manifesto for agile software development," Snowbird UT, 2001, [Online]. Available: https://agilemanifesto.org/

[3] M. Fowler and J. Highsmith, "The agile manifesto," Softw. Dev., vol. 9, no. 8, pp. 28–35, 2001.

[4] J. Hesselberg, Unlocking agility: An insider's guide to agile enterprise transformation. Addison-Wesley Professional, 2018. Accessed: Oct. 28, 2023. [Online]. Available: https://books.google.co.uk/books?hl=en&lr=&id=nS9mDwAAQBAJ&o i=fnd&pg=PT23&dq=Hesselberg+Unlocking+Agility.&ots=b-CQeV1oK-&sig=uuoTcm5CmebNRWXKWEQzCZk3Gco

[5] D. Cohen, M. Lindvall, and P. Costa, "An introduction to agile methods.," Adv Comput, vol. 62, no. 03, pp. 1–66, 2004.

[6] V. Szalvay, "An introduction to agile software development," Danube Technol., vol. 3, 2004, Accessed: Oct. 31, 2023. [Online]. Available: https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=2efe4 d840631ebf026fede741e85195e36f8b134

[7] S. Sharma, D. Sarkar, and D. Gupta, "Agile processes and methodologies: A conceptual study," Int. J. Comput. Sci. Eng., vol. 4, no. 5, p. 892, 2012.

[8] D. Taibi, V. Lenarduzzi, C. Pahl, and A. Janes, "Microservices in agile software development: a workshop-based study into issues, advantages, and disadvantages," in Proceedings of the XP2017 Scientific Workshops, Cologne Germany: ACM, May 2017, pp. 1–5. doi: 10.1145/3120459.3120483.

[9] A. DŽANIĆ, A. Toroman, and A. DŽANIĆ, "AGILE SOFTWARE DEVELOPMENT: MODEL, METHODS, ADVANTAGES AND DISADVANTAGES.," Acta Tech. Corviniensis-Bull. Eng., vol. 15, no. 4, 2022, Accessed: Oct. 31, 2023. [Online]. Available: https://acta.fih.upt.ro/pdf/2022-4/ACTA-2022-4-15.pdf

[10] H. Edison, X. Wang, and K. Conboy, "Comparing methods for large-scale agile software development: A systematic literature review," IEEE Trans. Softw. Eng., vol. 48, no. 8, pp. 2709–2731, 2021.

[11] D. E. Strode, "Agile methods: a comparative analysis," in Proceedings of the 19th annual conference of the national advisory committee on computing qualifications, NACCQ, 2006, pp. 257–264. Accessed: Oct. 31, 2023. [Online]. Available: https://www.researchgate.net/profile/Diane-Strode/publication/228918891_Agile_methods_a_comparative_analysis/ links/00b4951c7fbc8b72b9000000/Agile-methods-a-comparative-analysis.pdf

[12] K. Conboy and N. Carroll, "Implementing large-scale agile frameworks: challenges and recommendations," IEEE Softw., vol. 36, no. 2, pp. 44–50, 2019.

[13] F. Almeida and E. Espinheira, "Large-scale agile frameworks: a comparative review," J. Appl. Sci. Manag. Eng. Technol., vol. 2, no. 1, pp. 16–29, 2021.

[14] K. Schwaber and J. Sutherland, "The scrum guide," Scrum Alliance, vol. 21, no. 1, pp. 1–38, 2011.

[15] K. Schwaber and J. Sutherland, "The Scrum Guide. 2020," Accessed April, 2021, Accessed: Oct. 31, 2023. [Online]. Available:

https://topasspmp.com/wp-content/uploads/2021/01/SCRUM-GUIDE-2020-VIETNAMESE.pdf

[16] K. S. Rubin, Essential Scrum: A practical guide to the most popular Agile process. Addison-Wesley, 2012. Accessed: Oct. 31, 2023. [Online]. Available: https://books.google.co.uk/books?hl=en&lr=&id=3vGEcOfCkdwC&oi=fnd&pg=PR11&dq=%E2%80%9CThe+Scrum+GuideTM&ots=-DFesmcq1t&sig=aHzbVXXv0_BGuy8npnyZ3AQ0pFA

[17] K. Beck, "Embracing change with extreme programming," Computer, vol. 32, no. 10, pp. 70–77, 1999.

[18] K. Beck, Extreme programming explained: embrace change. addison-wesley professional, 2000. Accessed: Oct. 31, 2023. [Online]. Available: https://books.google.co.uk/books?hl=en&lr=&id=G8EL4H4vf7UC&oi=fnd&pg=PR13&dq=extreme+programming&ots=jbBLzqiStr&sig=wi8mXrHOrEc4k9oKHkUEGBfAnfc

[19] K. Beck and M. Fowler, Planning extreme programming. Addison-Wesley Professional, 2001. Accessed: Oct. 31, 2023. [Online]. Available: https://books.google.co.uk/books?hl=en&lr=&id=u13hVoYVZa8C&oi=fnd&pg=PR11&dq=eXtreme+Programming+&ots=GN5c1ScQdd&sig=vSQ3Vw1Qvc8NqurW3Kj0pgBnPC8

[20] G. Bou Ghantous and A. Gill, "DevOps: Concepts, practices, tools, benefits and challenges," PACIS2017, 2017, Accessed: Oct. 31, 2023. [Online]. Available: https://opus.lib.uts.edu.au/bitstream/10453/130066/1/DevOps-%20Concepts%20Practices%20Tools%20Benefits%20and%20Challenges.pdf

[21] M. Gall and F. Pigni, "Taking DevOps mainstream: a critical review and conceptual framework," Eur. J. Inf. Syst., vol. 31, no. 5, pp. 548–567, Sep. 2022, doi: 10.1080/0960085X.2021.1997100.

[22] W. P. Luz, G. Pinto, and R. Bonifácio, "Adopting DevOps in the real world: A theory, a model, and a case study," J. Syst. Softw., vol. 157, p. 110384, 2019.

[23] T. Dyba and T. Dingsoyr, "What do we know about agile software development?," IEEE Softw., vol. 26, no. 5, pp. 6–9, 2009.

[24] N. Ozkan, M. Ş. Gök, and B. Ö. Köse, "Towards a better understanding of agile mindset by using principles of agile methods," in 2020 15th Conference on Computer Science and Information Systems (FedCSIS), IEEE, 2020, pp. 721–730. Accessed: Oct. 28, 2023. [Online]. Available: https://ieeexplore.ieee.org/abstract/document/9222861/

[25] A. Przybyłek, M. Albecka, O. Springer, and W. Kowalski, "Game-based Sprint retrospectives: multiple action research," Empir. Softw. Eng., vol. 27, no. 1, p. 1, Oct. 2021, doi: 10.1007/s10664-021-10043-z.

[26] I. Nurdiani, J. Börstler, S. Fricker, K. Petersen, and P. Chatzipetrou, "Understanding the order of agile practice introduction: Comparing agile maturity models and practitioners' experience," J. Syst. Softw., vol. 156, pp. 1–20, 2019.

[27] H. Bundtzen and G. Hinrichs, "The link between organizational agility and VUCA–an agile assessment model," 2021, Accessed: Oct. 28, 2023.

[Online]. Available: https://essuir.sumdu.edu.ua/handle/123456789/83934

[28] G. Y. Koi-Akrofi, J. Koi-Akrofi, and H. A. Matey, "Understanding the characteristics, benefits and challenges of agile it project management: A literature based perspective," Int. J. Softw. Eng. Appl. IJSEA, vol. 10, no. 5, pp. 25–44, 2019.

[29] A. Hess, P. Diebold, and N. Seyff, "Understanding information needs of agile teams to improve requirements communication," J. Ind. Inf. Integr., vol. 14, pp. 3–15, 2019.

[30] T. Peeters, K. Van De Voorde, and J. Paauwe, "The effects of working agile on team performance and engagement," Team Perform. Manag. Int. J., vol. 28, no. 1/2, pp. 61–78, Jan. 2022, doi: 10.1108/TPM-07-2021-0049.

[31] L. Barroca, T. Dingsøyr, and M. Mikalsen, "Agile Transformation: A Summary and Research Agenda from the First International Workshop," in Agile Processes in Software Engineering and Extreme Programming – Workshops, R. Hoda, Ed., in Lecture Notes in Business Information Processing. Cham: Springer International Publishing, 2019, pp. 3–9. doi: 10.1007/978-3-030-30126-2_1.

[32] P. Kokol, "Agile Software Development in Healthcare: A Synthetic Scoping Review," Appl. Sci., vol. 12, no. 19, p. 9462, 2022.

[33] Ö. Uludağ, P. Philipp, A. Putta, M. Paasivaara, C. Lassenius, and F. Matthes, "Revealing the state of the art of large-scale agile development research: A systematic mapping study," J. Syst. Softw., vol. 194, p. 111473, 2022.

[34] J. E. Ravn, N. B. Moe, V. Stray, and E. A. Seim, "Team autonomy and digital transformation: Disruptions and adjustments in a well-established organizational principle," AI Soc., vol. 37, no. 2, pp. 701–710, Jun. 2022, doi: 10.1007/s00146-022-01406-1.

[35] J. Jia, X. Yang, R. Zhang, and X. Liu, "Understanding software developers' cognition in agile requirements engineering," Sci. Comput. Program., vol. 178, pp. 1–19, Jun. 2019, doi: 10.1016/j.scico.2019.03.005.

[36] A. R. Amna and G. Poels, "Systematic literature mapping of user story research," IEEE Access, vol. 10, pp. 51723–51746, 2022.

[37] C. Baham and R. Hirschheim, "Issues, challenges, and a proposed theoretical core of agile software development research," Inf. Syst. J., vol. 32, no. 1, pp. 103–129, 2022, doi: 10.1111/isj.12336.

[38] K. Eilers, B. Simmert, and C. Peters, "Doing agile vs. being agile-understanding their effects to improve agile work," 2020, Accessed: Oct. 28, 2023. [Online]. Available: https://www.alexandria.unisg.ch/server/api/core/bitstreams/fe287039-5e48-4aa9-8472-96711499b146/content

[39] D. Rad and G. Rad, "Going agile, a post-pandemic universal work paradigm-a theoretical narrative review," Postmod. Open., vol. 12, no. 4, pp. 337–388, 2021.

[40] J. Koch, I. Drazic, and C. C. Schermuly, "The affective, behavioural and cognitive outcomes of agile project management: A preliminary meta‐analysis," J. Occup. Organ. Psychol., vol. 96, no. 3, pp. 678–706, Sep. 2023, doi: 10.1111/joop.12429.

# Category Decomposition-based Within Pixel Information Retrieval Method and its Application to Partial Cloud Extraction from Satellite Imagery Pixels

Kohei Arai[1], Yasunori Terayama[2], Masao Moriyama[3]

Faculty of Science and Engineering, Saga University, Saga City, Japan[1, 2]

Faculty of Engineering, Nagasaki University, Nagasaki City, Japan[3]

*Abstract*—Category decomposition-based within pixel information retrieval method is proposed together with its application to partial cloud extraction from satellite imagery pixels. A comparative study was conducted for estimation of the sea surface temperature of the pixel suffered from partial cloud cover within a pixel. Three methods for estimation of partial cloud cover within a pixel, based on the proposed category decomposition-based method with Generalized Inverse Matrix Method: GIMM and well-known Least Square Method: LSM and Maximum Likelihood Method: MLH, were compared. It was found that around 9% of RMS (Root Mean Square) error can be achieved. Also, it was found that estimation accuracy highly depends on variance of representative vectors for cloud and the ocean or observed noise. The experimental results with simulated data show RMS error of GIMM are highly dependent to the noise followed by MLH and LSM. The results also show the best estimation accuracy can be achieved for MLH followed by LSM and GIMM.

*Keywords*—*Category decomposition; information retrieval; cloud cover estimation; Generalized Inverse Matrix Method: GIMM and well-known Least Square Method: LSM and Maximum Likelihood Method: MLH*

## I. INTRODUCTION

When estimating the sea surface temperature using visible thermal infrared radiometer data such as NOAA (National Oceanic and Atmospheric Administration) / AVHRR (Advance Very High-Resolution Radiometer), MOS-1 (Marine Observation Satellie-1) / VTIR (Visible and Thermal Infrared Radiometer), for example, as is clear from the MCSST (Multi-Channel Sea Surface Temperature) [1] algorithm, there is only a small amount in the pixel. However, the pixels that are likely to have clouds are detected and excluded from the target of sea surface temperature estimation. Especially in the case of MCSST, the acquisition rate of data not covered by clouds is low because the policy of punishing suspicions is strictly checked for this possibility. As a result, many observation day data are required to obtain a good scene in which all pixels are not covered with clouds, which often hinders the estimation of the 10-day average sea surface temperature.

Even if a small cloud exists in the pixel, if the brightness temperature of the cloud can be known and the area occupancy can be estimated, it can be corrected to some extent and used. Assuming that the brightness temperature of this cloud is equal to that of the pixel covered with 100% cloud in the vicinity of the core pixel, the method of estimating the cloud coverage rate will be examined here. That is, with the aim of creating products of average sea surface temperature in a short period of time, we propose a method for estimating the cloud coverage rate in pixels and examine its effect.

In this paper, we take up the method of estimating the class occupancy in pixels proposed so far as a method of estimating the cloud coverage [2]-[8] and show the result of mutual comparison of estimation accuracy. These estimation methods have been proposed to estimate the class mixing ratio of mixed pixels (Mixel) consisting of multiple classes. When applying these to cloud coverage estimation, it becomes a problem to estimate the mixing ratio for the two classes of cloud and sea, and in general, the number of channels of visible thermal infrared radiometer data exceeds this, so the minimum square method is effective. It is considered to work. Therefore, we took up the least squares method that minimizes the square of the estimation error of the observation vector and the square of the estimation error of the mixing ratio. In addition, we conducted a theoretical study of the estimation error of these least square methods, and the estimation error is small.

It is shown that it is possible to use both adaptively so as to become. We propose an "adaptive least squares method" based on that principle and apply the effect to actual data to confirm it. Furthermore, since it is expected that the spectral reflection and radiation characteristics of cloud pixels will vary widely, the maximum likelihood method that takes the variance into consideration was taken up as a comparison target and compared.

In the next section, related research works are described in Section II followed by theoretical background and proposed method in Section III. Experiments and experiments results are mentioned in Section IV and Section V respectively and finally conclusion and work for future is explained in Section VI and Section VII respectively.

## II. RELATED RESEARCH WORKS

As for the related research works to category decomposition, there are the followings,

Maximum likelihood estimation of category proportion among Mixels is conducted [9]. Meanwhile, image classification from category proportions among Mixels is proposed [10]. On the other hand, decomposition of category

mixture in a pixel and its application for supervised image classification is proposed [11].

Category decomposition based on subspace method with learning process is proposed [12] together with category decomposition method for un-mixing of Mixels acquired with spaceborne based visible and near infrared radiometers by means of Maximum Entropy Method: MEM with parameter estimation based on Simulated Annealing: SA [13]. On the other hand, focusing on Mixels located at the boundary between two types of classes, an image decomposition algorithm that uses the class mixture ratio of Mixels and the spatial information of surrounding pixels of Mixels are investigated [14]. Research has also been conducted that adds, but it has not been applied to two or more types of class boundaries. Meanwhile, category decomposition requires an endmember extraction in the spectral space of distributions. When observation data, end member spectra, and content rates are each expressed as a matrix, Mixel decomposition can be regarded as a matrix decomposition problem. Due to physical conditions, all components of the matrix are non-negative values, so by applying non-negative matrix factorization (NMF), the end member spectrum and content can be estimated simultaneously [15]. In the data-driven approach, the material with the estimated endmember spectrum is finally identified by referring to the spectral library and searching for the closest spectrum.

Sea ice concentration estimation method with satellite based visible to near infrared radiometer data based on category decomposition is proposed [16]. Also, category decomposition method based on matched filter for un-mixing of mixed pixels acquired with space borne based hyper-spectral radiometers is proposed [17].

Bi-directional Reflectance Distribution Function: BRDF effect on un-mixing, category decomposition of the Mixel of remote sensing satellite imagery data is estimated [18].

On the other hand, there are the following research works related to cloud overage estimation,

A merged dataset for obtaining cloud free Infrared: IR data and a cloud cover estimation within a pixel for SST retrieval is proposed [19]. Meanwhile, estimation of partial cloud coverage within a pixel is conducted [20].

Comparative study on estimation of partial cloud coverage within a pixel is conducted [21]. On the other hand, adjacency effect of layered clouds estimated with Monte-Carlo simulation is estimated [22].

Evaluation of cirrus cloud detection accuracy of GOSAT/CAI (Green House Gasses Observation Satellite / Cloud and Aerosol Imager) and Landsat-8 with laser radar: lidar and confirmation with CALIPSO (Cloud-Aerosol Lidar and Infrared Pathfinder Satellite Observations) data is conducted [23]. Meanwhile, comparative study on cloud parameter estimation among GOSAT/CAI, MODIS (Moderate Resolution Imaging SpectroRadiometer), CALIPSO/CALIOP (Cloud-Aerosol LIdar with Orthogonal Polarization) and Landsat-8/OLI (Operational Land Imager - Landsat Science) with laser radar as truth data is conducted [24].

Thresholding-based method for rain, cloud detection with NOAA/AVHRR data by means of Jacobi iteration method is proposed [25]. Also, adjacency effects of layered clouds by means of Monte Carlo Ray Tracing: MCRT is investigated [26].

## III. THEORETICAL BACKGROUND AND PROPOSED METHOD

### A. Category Decomposition and Classification Norms

In order to estimate the maximum occupancy category in Mixel, category decomposition [1] is required to estimate the occupancy rate of each category. Several categorical decomposition methods have been devised [1-6], but in this study, the categorical decomposition is formulated using the maximum likelihood estimation method that takes observation errors into consideration. Using this theory has the advantage that unclassified pixels can be determined to be statistically meaningful.

Mixel's spectroscopic vector: *I*, which is a mixture of information from N categories, is considered to be the linear combination of Pure pixel value: *A* and category occupancy: *B* shown in Eq. (1) plus the observation error vector: ε.

$$I = (I_1, I_2, \dots, I_M)^t = AB + \varepsilon \tag{1}$$

where *Ii*: Observation pixel value of the i-th band, *M*: Number of bands. "t" represents transpose. Furthermore, *A* is expressed as follows:

$$A = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & \ddots & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix}$$

where *Aij*: i-band of pure pixel value of j-category, $B = (B_1, B_2, \cdots, B_N)^t$, *Bj*: Occupancy of category j. And, $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_M)^t$, εi: observation error of the i-th band. Here, it is assumed that Aij follows the normal distribution of mean $A^*_{ij}$ and variance $\sigma ij^2$: N $(A^*_{ij}, \sigma_{ij}^2)$, and $\varepsilon_i$ follows N $(0, \sigma_{ei}^2)$, and the spectroscopic vector 1 is a random variable [4]. Here, if the pixel values of Pure pixels in each category are independent, the observed pixel values of the i-band: *I*, are the average $A^*_i$ represented by Eq. (2) and Eq. (3), and the normal distribution of the variance $\sigma_i^2$: N $(A^*_i, \sigma_i^2)$ is obeyed [7].

$$I_i = A_i^* B \tag{2}$$
$$A_i^* = (A_{i1}^*, \dots, A_{iN}^*)$$
$$\sigma_i^2 = B^t S_i B + \sigma_{ei}^2 \tag{3}$$
$$S_i = diag(\sigma_{i1}^2, \dots, \sigma_{iN}^2)$$

where $A^*_{ij}$: The average of the pure pixel values of the i-band j category, and $\sigma_{ij}^2$: the variance of the pure pixel values of the i band, j-category. It is also the variance of the observation error of the $\sigma_{ei}^2$: i band.

Observed pixel value of the i-th band: Probability (likelihood) that *Ii* is observed: *P (Ii)* is expressed by Eq. (4).

$$P(I_i) = \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(I_i - A_i^*)^2}{2\sigma_i^2}} \tag{4}$$

Probability (likelihood) that the spectroscopic vector *I* is observed (likelihood): *P(I)* is expressed by Eq. (5), assuming that it is independent between each band.

$$P(I) = \prod_{i=1}^{M} P(I_i) \qquad (5)$$

In general, each band of remote sensing data is not independent of each other, so it is necessary to make each band independent by orthogonalization transformation such as principal component analysis as preprocessing. Here, based on the concept of the maximum likelihood estimation method, the solution is to find the probability that the spectral vector is observed: *P(I)* and the occupancy rate: *B I*. Here, the occupancy rate: *B* has the following constraints if the type of category included in Mixel is known.

$$\sum_{j=1}^{N} B_j = 1, B_j \geq 0, (j = 1, \dots, N) \qquad (6)$$

The category occupancy (maximum likelihood estimation value) estimated based on this method is expressed as $B^*$, and the pixels are classified into category *K* of the maximum element $B^*_k$ in $B^*$. This method is called the Maximum Proportion Classifier (MPC).

### B. Determination of Unclassified Pixels

In image classification methods such as the maximum likelihood method and the shortest distance method, restrictions are set according to those classification norms, and those exceeding the restrictions are regarded as unclassified pixels [8], [9]. This section describes how to determine unclassified pixels according to the maximum occupancy classification norm.

Comparing Mixel and Pure Pixel, Pure Pixel expresses its pixel value by the mean and variance of the classified categories, while Mixel's pixel value is the occupancy rate of each category and their mean and variance. This indicates that Mixel has more independent parameters to express the pixel value of Mixel than Pure Pixel, and Mixel is a concept with a high degree of freedom. With these things in mind, this study proposes the following method for determining unclassified pixels.

The fact that a Mixel can be classified into one category (Mixel can be represented by one category) means that the Mixel can be represented by a model with more constraints (fewer independent parameters) (Pure pixel hypothesis). That's what it means. In such a case, the Pure pixel hypothesis is tested using the model goodness-of-fit test [7], [10], and if the Pure pixel hypothesis holds, it can be classified, otherwise it cannot be classified. It is possible to do. Here, we propose two types of goodness-of-fit test methods.

### C. Goodness of Fit Determined by $\chi^2$ Distribution [10]

Suppose there are two models $\pi_1$ and $\pi_2$, and $\pi_1$ is $\pi_2$, which is a special case (the number of independent parameters is small). The likelihood ratio of the two models: β is defined by Eq. (7), where P (π) is the likelihood of the model π.

$$a = P (\pi^*_1) / P (\pi^*_2) \qquad (7)$$

where, the superscript * represents the maximum likelihood estimator of the model $\pi$. Since $\pi^*_1$ is a special case of $\pi^*_2$,

P ($\pi^*_1 \leq$ P ($\pi^*_2$)), and therefore $\beta \leq 1$ holds. Here, $\chi^2$ of Eq. (8) is defined.

$$\chi^2 = -2 \ln \beta \qquad (8)$$

$\chi^2$ asymptotically has a chi-square distribution with $n = n_2 - n_1$ degrees of freedom. Here ($n_1$ and $n_2$ are the number of independent parameters of the models $\pi_1$ and $\pi$, respectively), and the percentile value $\chi^2$ (n) of 100 by α% (α: significance level, $0 < \alpha < 1$) of the chi-square distribution with n degrees of freedom, $\alpha$) and $\chi^2$ can be compared to test whether π, is significantly inferior to $\pi_2$ ($\chi^2$ (n, α) $< \chi^2$).

When this goodness-of-fit test is used to determine unclassified pixels according to the maximum occupancy classification standard, only the maximum occupancy category $B^*_k$ of the most likely estimated value $B^*$ of the occupancy obtained by categorization of $\pi_1$ Pire pixel ($B^*_K = 1$, occupancy of other elements is 0), $P(\pi_1)$ is obtained from Eq. (4) and Eq. (5), then $\pi_2$ is Mixel, and B* is used (4). ), (5) to find P ($\pi_2$), and Eq. (7) and Eq. (8) to find $\chi^2$. Here, from the constraint condition of Eq. (6), the number of independent parameters of PURE PIXEL is 0, and the number of independent parameters of Mixel is *N*-1 (*N*: number of categories), so n=N-1. Therefore, the significance level α (right side test) of the test is determined, $\chi^2$ (*N*-1,α) is calculated, compared with $\chi^2$, and the unclassified pixels are $\chi^2$ (N-1),α) $\leq \chi^2$: Unclassified $\chi^2$ (N-1,α)> $\chi^2$: Determined to be classified in the maximum occupancy category. Since $\chi^2$ becomes larger as the likelihood ratio $\beta$ is smaller (the likelihood is smaller when it is a pure pixel), the number of unclassified pixels increases as the significance level is increased ($\chi^2$ (*N*-1, α) becomes smaller).

### D. Goodness of Fit Test by AIC [10],[11]

Similar to the goodness-of-fit test based on the $\chi^2$ distribution, there are two models, $\pi_1$ and $\pi$, and $\pi_1$ is a special case of $\pi_2$. If $P(\pi)$ is the likelihood of the model π, then AIC (Akaike's Information Criterion) is defined by Eq. (9).

$$AIC = 2(n - 2)\ln\{P(\pi^*)\} \qquad (9)$$

where *n* is the number of independent parameters of the model, and the superscript * is the maximum likelihood estimator of the model. Here, the model that minimizes the AIC in Eq. (9) is selected. Similar to the goodness-of-fit test based on the $\chi^2$ distribution, $\pi_2$ is Mixel (number of independent parameters: *N*-1), $\pi_1$ is Pure pixel (number of independent parameters: 0) containing only the maximum occupancy category, and AIC is Eq. (4), (5), (9), if the AIC obtained from the case of Pure pixel is smaller than the AIC obtained from the case of Mixel, it is judged that it can be classified and classified into the maximum occupancy category. If not, it is regarded as an unclassified pixel.

## IV. EXPERIMENTS

### A. Simulation of Pure Pixel Data

The validity of the above theory will be confirmed by the following data and simulation according to the procedure.

Using TM data around Lake Ashino-ko in Japan, which was acquired by LANDSAT 5 on June 6, 1987, residential areas, bare land, grasslands, coniferous forests, and broad-

leaved forests, which are typical categories of this area, were ed as the grand truth area. It is extracted using vegetation maps and aerial photographs [12]. Then, in order to make each band independent, all bands except the thermal band (band 6) were analyzed for principal components, and the first and second principal components were quantized into eight bits and used for categorization. Table I shows the eigenvalues, eigenvectors, and contribution ratio of each principal component. Table II shows the average and variance of the pixel values of each category Pure pixel.

TABLE I.    EIGENVALUE, EIGENVECTOR AND CONTRIBUTION OF THE TEST DATA. PC, E-VALUE, E-VEC. AND CONT. MEAN PRINCIPAL COMPONENT, EIGENVALUE, EIGENVECTOR AND CONTRIBUTION, RESPECTIVELY

|  | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 |
|---|---|---|---|---|---|---|
| E-value | 5.47 | 0.49 | 0.03 | 0.01 | 0.00 | 0.00 |
| Band 1 | 0.41 | -0.33 | 0.00 | -0.83 | 0.13 | -0.08 |
| Band 2 | 0.42 | -0.26 | -0.27 | 0.32 | -0.35 | -0.68 |
| Band 3 | 0.42 | -0.29 | -0.12 | 0.19 | -0.42 | 0.72 |
| Band 4 | 0.36 | 0.75 | -0.53 | -0.12 | 0.07 | 0.06 |
| Band 5 | 0.41 | 0.39 | 0.79 | 0.02 | -0.23 | -0.09 |
| Band 7 | 0.42 | -0.16 | 0.09 | 0.39 | 0.79 | 0.06 |
| Cont. | 0.91 | 0.08 | 0.01 | 0.00 | 0.00 | 0.00 |

TABLE II.    AVERAGE AND VARIANCE OF THE PURE PIXEL DATA. AV AND VR MEANS AVERAGE AND VARIANCE, RESPECTIVELY. THE NUMBERS 1 TO 5 SHOW THAT THE CATEGORIES OF RESIDENTIAL AREA, BARE SOIL, GRASS LAND, NEEDLE LEAF TREE AND BROAD LEAF TREE, RESPECTIVELY

|  |  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| A | PC1 | 97.8 | 162.4 | 127.3 | 60.9 | 107.8 |
| V | PC2 | 62.2 | 135.1 | 162.0 | 100.9 | 187.7 |
| V | PC1 | 160.4 | 841.1 | 185.7 | 94.0 | 178.2 |
| R | PC2 | 309.9 | 681.3 | 430.4 | 329.3 | 586.2 |

### B. Simulation of Mixel Data

Mixel data of 36 types (six types of variances by six types of observation error) were created by the procedure shown below,

*1) Category occupancy rate:* Using the uniform random numbers of [0,1), 100 points of category occupancy are created based on the constraint condition of Eq. (6).



(a)                                     (b)

Fig. 1.   Characteristics of the truth data. (a) Histogram of the maximum proportion category. The categories number 1 to 5 correspond to the categories of Residential area, Bare soil, Grass land, Needle leaf tree and Broad leaf tree, respectively. (b) Histogram of the maximum proportion.

Fig. 1(a) and Fig. 1(b) show the distribution of the maximum occupancy category and the distribution of the maximum occupancy, respectively.

*2) Pure pixel value (training data):* Pure pixel values according to the mean and variance of each category are created for each band using normal random numbers. At this time, in order to confirm the influence of the variance of the Pure pixel value, the Pure pixel value is created by multiplying the variance of each category and each band by 2.0, 1.0, 0.8, 0.6, 0, 4, 0.2, respectively.

*3) Observation error:* The standard deviation $\sigma_{ei}$ of the observation error in each band is 0,2,4,6,8,10 [Count]. The observation error is created using normal random numbers.

*4) Mixel dataset:* From the data of (1)-(3), 100 points of Mixel spectroscopic vectors are created under each condition according to Eq. (1).

### C. Verification Details

The 30 types of Mixel data sets created by the above method were categorized using the grid search method [13] with a side length of 1/64. This solution creates a mesh with a side length of 1/64 in the solution space ($N$-1 dimensional hyperplane) given by Eq. (6) and uses the training data given at each point of the mesh and the observed spectral vector. The likelihood given by Eq. (5) is calculated, and the point that gives the maximum likelihood is the solution.

As for the solution of the nonlinear optimization problem, all high-speed calculation methods such as Newton's method are methods for finding extreme values, not methods for finding maximum / minimum values. In this study, in order to avoid a decrease in the accuracy of categorization due to algorithm restrictions, categorization was performed by the lattice search method without using high-speed calculation. In actual applications, it is necessary to develop accurate and high-speed algorithms, which will be an issue for the future.

From the estimated value of the category occupancy and the estimated maximum occupancy category obtained here, the following items are verified together with the truth data of the category occupancy in Eq. (1) Category occupancy rate.

### D. Comparison of Goodness-of-Fit Test by $\chi^2$ Distribution and Goodness-of-Fit Test by AIC

In order to compare the two unclassified pixel determination methods and confirm the effect of the goodness-of-fit test based on the $\chi^2$ distribution on the classification accuracy of water a, the goodness-of-fit test by AIC and the significance level α were set to 1, 5, 10%. Classification was performed using MPC to which the method for determining unclassified pixels by the goodness-of-fit test was applied. In this case, the same method for determining unclassified pixels is applied to the truth data of the occupancy rate of Eq. (1) Category occupancy rate, and only the classifiable pixels are extracted, and the classifiable pixels obtained here, and their maximum occupancy category are classified and was used as the truth data of the category to be classified.

*E. Comparison with Maximum Likelihood Method*

The 36 types of Mixel datasets created in steps 3 and 2 were classified by the maximum likelihood classifier (MLC) and compared with the classification results by MPC. In this case, the following two types of unclassified pixel determination methods were used. We used to extract classifiable pixels.

*1)* Not classified if the log-likelihood is -20 or less.

*2)* If the Mahalanobis distance to the target category is three times or more the maximum standard deviation of that category, it is unclassified.

Here, as in Eq. (1) Category occupancy rate, the true data of the occupancy rate is classified by MPC using the unclassified pixel determination method by AIC, and the classifiable pixels and their maximum occupancy rate categories are extracted and classified and the truth data of the category to be used.

## V. Experimental Results

*A. Comparison of Goodness-of-Fit Test by $\chi^2$ Distribution and Goodness-of-Fit Test by AIC*

The occupancy rate estimated by categorical decomposition and the likelihood used to determine unclassified pixels are affected by the variance term $\sigma_1$ in Eq. (3). Since the variance term is determined by the category occupancy rate, the variance of the Pure pixel value, and the variance of the observation error, the average variance AVG [$\sigma$] of Eq. (10) is obtained for each of the 30 types of Mixel data (100 points each) and used. The characteristics of each Mixel dataset are shown.

$$AVG[\sigma] = \frac{1}{100}\sum_{k=1}^{100}\{\frac{1}{M}\sum_{i=1}^{M}(B_k^t S_i B_k + \sigma_{ei}^2)\} \quad (10)$$

First, in order to confirm the accuracy of categorization, the root mean square error RMSE between the occupancy rate $B^*$ and the true value $B$ estimated from each Mixel data set is obtained from Eq. (11), and the mean variance AVG [$\sigma$] is used. The relationship is shown in Fig. 2.



Fig. 2. The relationship between the mean variance and the root mean square error of the estimated category proportions from maximum likelihood estimation (MLE) and generalized inversion matrix (GIM).

Here, the results when the general inverse matrix, which is a typical conventional categorical decomposition method, is used are also shown.

$$RMSE = \frac{1}{100}\sum_{k=1}^{100}\{\frac{1}{N}(B_k^* - B_k)^t(B_k^* - B_k)\} \quad (11)$$

In both methods, it was confirmed that as the mean variance increases (the variance of the Pure pixel value and the observation error variance increase), the RMSE increases and the accuracy of categorization decreases. In addition, the generalized inverse matrix does not take into account the variation in the Pure pixel values of each category [1], indicating that the estimation accuracy is lower than that of the maximum likelihood estimation method.

Next, in order to clarify the relationship between the goodness-of-fit test method for determining unclassified pixels and the degree to which the Mixel data set created by the above method is judged to be classifiable, these unclassified truth data are included in the truth data of the occupancy rate. The relationship between the number of classable pixels obtained by applying the classification pixel determination method and the average variance was obtained. The results are shown in Fig. 3.



Fig. 3. The relationship between the mean variance and the number of selected classifiable pixel from the truth data of the proportion from Chi square and AIC method.

As the mean variance increases, the likelihood P represented by Eq. (4) and Eq. (5) becomes a gentle function, so the number of pixels judged to be classifiable by all goodness-of-fit test methods increases. In addition, it was confirmed that the number of pixels judged to be classable decreases as the significance level increases in the goodness-of-fit test using the $\chi^2$ distribution, and the goodness-of-fit test using the AIC is a goodness-of-fit test using the $\chi^2$ distribution with the significance level set to 10%.

It was confirmed that almost the same result as the above was obtained. Here, as classification accuracy verification, from the classification result using the estimated value of the category occupancy rate, it was judged that (a) the number of pixels judged to be categorizable and (b) the truth data of the occupancy rate could be classified. The number of pixels

determined to be unclassified by the estimated value (number of error pixels of the first type), (c) Pixels determined to be unclassified by the truth data of the occupancy rate but determined to be classifiable by the estimated value The number (the number of pixels of the type II error) and (d) the classification correctness: (the number of pixels correctly classified) / (the number of pixels judged to be classifiable) were calculated. The results are shown in Fig. 4(a) to Fig. 4(d).



Fig. 4. The classification results from Maximum proportion classifier from various unclassified limits. (a) The number of selected classifiable pixels. (b) The number of unselected classifiable pixels. (c) The number of mis-selected classifiable pixels. (d) The ratio of correctly classified pixels.

It was confirmed that the number of pixels judged to be classable increased as the mean variance increased and decreased as the significance level increased, as in the case obtained from the truth data of the occupancy rate, and the goodness-of-fit test by AIC was significant. It was confirmed that almost the same result as the goodness-of-fit test based on the $\chi^2$ distribution when the level was set to 10% was obtained.

It was confirmed that the number of pixels of the type I error increases with the increase of the average variance up to about 100, and then converges or decreases comparatively gently. In the goodness-of-fit test using the $\chi^2$ distribution, when the significance level was reduced, the number of pixels of type I errors tended to decrease rapidly when the mean variance increased. This can be explained by the fact that the number of pixels judged to be categorizable increases as the significance level increases. In this case as well, the goodness-of-fit test by AIC gave almost the same results as the goodness-of-fit test by $\chi^2$ distribution when the significance level was 10%.

It was confirmed that the number of pixels of the type II error increases as the average variance increases. It was also confirmed that when the significance level was reduced in the goodness-of-fit test using the $\chi^2$ distribution, the number of pixels of the type II error decreased when the mean variance increased. This is because if the significance level is reduced, it is judged that it can be almost classified even when applied to the truth data of the occupancy rate. In this case as well, the goodness-of-fit test by AIC gave almost the same results as the goodness-of-fit test by $\chi^2$ distribution when the significance level was 10%.

The classification correctness tended to decrease as the mean variance increased, and it was confirmed that the difference in the correctness due to the difference in the goodness-of-fit test method also decreased as the mean variance increased. In this case, if the number of pixels determined to be classable is 0, the correctness rate is set to 0.

### B. Comparison with Maximum Likelihood Method

The result of classifying the above 30 Mixel datasets by the maximum likelihood method (MLC) with two unclassified limits (log-likelihood and Mahalanobis distance) and the number of unclassified pixels using the AIC goodness-of-fit test. A comparison of the classification results by MPC to which the determination method is applied is shown.

As classification accuracy verification, from the classification result using the estimated value of the category occupancy rate, (a) the number of pixels judged to be categorizable (b) the truth data of the occupancy rate was judged to be categorizable, but the estimated value is not yet.



Fig. 5. Comparison of classification result from Maximum proportion classifier with AIC based unclassified limit: MPC(AIC), Maximum likelihood classifier with likelihood based unclassified limit: MLS(L) and with distance based unclassified limit: MLC(D). (a) The number of selected classifiable pixels. (b) The number of unselected classifiable pixels. (c) The number of mis-selected classifiable pixels. (d) The ratio of correctly classified pixels.

The number of pixels determined to be classified (type I error pixels), (c) The number of pixels determined to be unclassified in the truth data of the occupancy rate but determined to be classified by the estimated value (second) (Number of pixels with type error), (d) Classification correctness: (Number of pixels correctly classified) / (Number of pixels judged to be classifiable) were calculated. The results are shown in Fig. 5(a) to Fig. 5(d).

It was confirmed that the number of pixels judged to be classable by the maximum likelihood method applying the unclassified pixel determination method using the likelihood increased with the increase of the average variance, but the average variance was about 150, and all of them. It is judged that it can be classified.

Fig. 5(a) shows that when the maximum likelihood method is applied to Mixel, which targets categories with large variance, it is difficult to set the unclassified limit. The effect of the unclassified limit setting value appears in the error analysis of the above, especially in the number of pixels of the type I and type II errors.

The number of error pixels of the first type (the number of pixels that are determined to be unclassified pixels even though they can be dispersed) is 0 because the number of pixels that are determined to be classifiable in the most probable method is large. The number of error pixels (the number of pixels judged to be classifiable even though they are unclassified pixels) tends to increase when the average variance is small (200 or less) compared to when MPC is used. This indicates that the maximum likelihood method increases misclassification when the mean variance is small (the variance of the Pure pixel value is large, and the variance of the observation error is large) when compared with the result of classification by MPC.

The classification correctness rate when the maximum likelihood method is used is about 30%, which is not so affected by the mean variance. On the other hand, the results by MPC show that the correctness rate decreases as the mean variance increases, but the correctness rate is generally higher than that by the maximum likelihood method. These results represent the limits of the maximum likelihood method, which assumes that the pixel is a pure pixel, and show the usefulness of the proposed method.

In addition, the unclassified limit in the maximum likelihood method does not mean that the likelihood, or the variance and classifiable of a particular category, is fully meaningful, and in addition, the parameters of unclassified pixel determination are for each category. Since it is sensitive to the dispersion of Pure pixel values, it is difficult to determine the optimum parameters.

On the other hand, the method of determining unclassified pixels by the goodness-of-fit test corresponding to MPC proposed in this study is based on the hypothesis test that the Mixel can be regarded as a pure pixel, and in addition, the pure pixel of each category. It is effective and easy to use because it is insensitive to pixel values. Furthermore, since the goodness-of-fit test by AIC is a method that excludes the arbitrariness of

the significance level, it has the advantage that unclassified pixels can be uniquely determined.

## VI. CONCLUSION

The following conclusions can be drawn from the above results. Considering that each pixel in remote sensing is a Mixel, we propose a maximum occupancy classification norm that classifies pixels into the maximum occupancy category, and in addition, a method for determining unclassified pixels based on the goodness of fit of the pixel as a pure pixel. It showed that from the simulation of artificially creating Mixel, the result that the proposed method has better classification accuracy than the maximum likelihood method was obtained, and the limit of the maximum likelihood method and the effectiveness of the proposed method were shown.

We also proposed two methods for determining unclassified pixels by the goodness-of-fit test according to the proposed method, one based on the $\chi^2$ distribution and the other based on the AIC, and it was confirmed that there was not much difference between the two. From this, it was concluded that the $\chi^2$ distribution should be used when the number of classified pixels should be adjusted according to the user's situation, and the AIC should be used when the significance level should be excluded.

Since this method takes a long time to calculate at present, it is used as a secondary application such as a remedy for pixels determined to be unclassified in the maximum likelihood classification, or classification of clouds and the sea in the sea area. It can be used as a classification method when the number of categories is small.

## VII. FUTURE RESEARCH WORKS

In the future, we plan to develop a high-speed calculation method for maximum likelihood estimation of category occupancy so that it can be used as a general classification method.

## REFERENCES

[1] Corniilon, P. et al ,. Sea Surf aco Tcmpcraturo Products for the Oceano graphic Scientific Research Community, Joint Oceanographic Institutions. Inc. p.1-36.1989.

[2] Inamura, 1987. Analysis of remote sensing image data based on category decomposition, Journal of the Institute of Electronics, Information and Communication Engineers, vo1.J70-C, No.2, pp.241-250 1987.

[3] Rikimaru, Uegami, Oshima, 1988. Development of a simple estimation method for intra-pixel spectroscopic information, Journal of the Japanese Society of Photometric Survey, vol.27, No.6, pp.23-34, 1988.

[4] Matsumoto, Fujiku, Tsuchiya, Arai, Category based on maximum likelihood estimation method, decomposition, Journal of Photogrammetry, Japan, Vol.30, No.2, pp.25-34,1991.

[5] Ito, Fujimura, Area Ratio Estimate by Pixel Category Decomposition, Proceedingseedings of the Society of Instrument and Control Engineers, Vol.23, No.8, pp.20-25,1987.

[6] Kohei Arai,and Y. Terayama. 1990. A Method for Proportion Estimate by Means of inversion Problem Solving, Proceedingseedings of the ISPRS Mid-Term Symposium, Commission Vil, WP-1-1, 1990.

[7] Kohei Arai, Terayama, Matsumoto, Fujiku, Tsuchiya, 1991, Context classification with class mixing ratio estimation in adjacent boundary pixels, Journal of the Remote Sensing Society of Japan, Vol.11, No.4, pp.21-28, 1991.

[8] Kohei Arai, Estimation of partial cloud coverage within a pixel, 1991. Proceedingseedings of the Pre-ISY International Symposium pp.97-106, 1991.

[9] M.Matsumoto, Y.Terayama and Kohei Arai, Maximum likelihood estimation of category proportion among mixels, Proceedings.of the ISPRS Committee-III, PS-6-8, 1992.

[10] M.Matsumoto, Y.Terayama and Kohei Arai, Image classification from category proportions among mixels, Proceedings.of the ISPRS Committee-VII, PS-9-15, 1992.

[11] M.Matsumoto, Kohei Arai, T.Ishimatsu, Decomposition of category mixture in a pixel and its application for supervised image classification, Proceedings.of the KACC'92, 514-519, 1992.

[12] Kohei Arai and Chen H., Category decomposition based on subspace method with learning proceedingsess, Abstract, COSPAR A1.1, A-00712, 2006.

[13] Kohei Arai, Category decomposition method for un-mixing of mixels acquired with spaceborne based visible and near infrared radiometers by means of Maximum Entropy Method with parameter estimation based on Simulated Annealing, International Journal of Advanced Research in Artificial Intelligence, 2, 4 64-69, 2013.

[14] Makoto Nishida, Yoichi Kageyama, Takashi Soma, Image Resolution Algorithm for Mixed Pixel of Boundary Area, T. IEE Japan, Vol. 116-C, No. 12, 1418-1419, 1996.

[15] Miao, L. and Qi, H.: Endmember extraction from highly mixed data using minimum volume constrained nonnegative matrix factorization, IEEE Trans. Geosci. Remote Sens., Vol. 45, No. 3, pp. 765-777, 2007.

[16] Kohei Arai, Sea ice concentration estimation method with satellite based visible to near infrared radiometer data based on category decomposition, International Journal of Advanced Research in Artificial Intelligence, 2, 5, 7-13, 2013.

[17] Kohei Arai, Category decomposition method based on matched filter for un-mixing of mixed pixels acquired with space borne based hyper-spectral radiometers, International Journal of Advanced Research in Artificial Intelligence, 2, 6, 20-26, 2013.

[18] Kohei Arai, Bi-directional reflectance distribution function: BRDF effect on un-mixing, category decomposition of the mixed pixel (MIXEL) of remote sensing satellite imagery data, International Journal of Advanced Research in Artificial Intelligence, 2, 9, 19-23, 2013.

[19] Kohei Arai, A Merged Dataset for Obtaining Cloud Free IR Data and a Cloud Cover Estimation within a Pixel for SST Retrieval, Asian-Pacific Remote Sensing Journal, Vol.4, No.2, pp.121-127, Jan.1992.

[20] Kohei Arai, Estimation of partial cloud coverage within a pixel, Proceedings.of the Pre-ISY International Symposium, 99-106, 1991.

[21] Kohei Arai, Y.Ueda and Y.Terayama, Comparative study on estimation of partial cloud coverage within a pixel, Proposed adaptive least square method with constraints- Proceedings.of the European ISY Conference, 305-310, 1992.

[22] Kohei Arai, Adjacency effect of layered clouds estimated with Monte-Carlo simulation, Advances in Space Research, Vol.29, No.19, 1807-1812, 2002.

[23] Kohei Arai, Masanori Sakashita, Evaluation of Cirrus Cloud Detection Accuracy of GOSAT/CAI and Landsat-8 with laser Radar: Lidar and Confirmation with Calipso Data, International Journal of Advanced Research on Artificial Intelligence, 5, 1, 14-21, 2016.

[24] Kohei Arai, Masanori Sakashita, Hiroshi okumura, Shuichi Kawakami, Kei Shiomi, Hirofumi Ohyama, Comparative Study on Cloud Parameter Estimation AmongGOSAT/CAI, MODIS, CALIPSO/CALIOP and Landsat-8/OLI with Laser Radar as Truth Data, International Journal of Advanced Research on Artificial Intelligence, 5, 5, 21-29, 2016.

[25] Kohei Arai, Thresholding Based Method for Rain, Cloud Detection with NOAA/AVHRR Data by Means of Jacobi Iteration Method, International Journal of Advanced Research on Artificial Intelligence, 5, 6, 21-27, 2016.

[26] Kohei Arai, Adjacency effects of layered clouds by means of Monte Carlo Ray Tracing, International Journal of Advanced Computer Science and Applications IJACSA, 11, 1, 95-98, 2020.

AUTHOR'S PROFILE

Kohei Arai, He received BS, MS and PhD degrees in 1972, 1974 and 1982, respectively. He was with The Institute for Industrial Science and Technology of the University of Tokyo from April 1974 to December 1978 also was with National Space Development Agency of Japan from January, 1979 to March, 1990. During from 1985 to 1987, he was with Canada Centre for Remote Sensing as a Post Doctoral Fellow of National Science and Engineering Research Council of Canada. He moved to Saga University as a Professor in Department of Information Science on April 1990. He was a councilor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology during from 1998 to 2000. He was a councilor of Saga University for 2002 and 2003. He also was an executive councilor for the Remote Sensing Society of Japan for 2003 to 2005. He is a Science Council of Japan Special Member since 2012. He is an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Science Commission "A" of ICSU/COSPAR during 2008 and 2020 then he is now award committee member of ICSU/COSPAR. He is now Visiting Professor of Nishi-Kyushu University since 2021, and is Visiting Professor of Kurume Institute of Technology (Applied AI Laboratory) since 2021. He wrote 87 books and published 700 journal papers as well as 570 conference papers. He received 66 of awards including ICSU/COSPAR Vikram Sarabhai Medal in 2016, and Science award of Ministry of Mister of Education of Japan in 2015. He is now Editor-in-Chief of IJACSA and IJISA. http://teagis.ip.is.saga-u.ac.jp/index.html

# Costless Expert Systems Development and Re-engineering

Manal Alsharidi, Abdelgaffar Hamed Ali

The Department of Information System, College of Computer Sciences and Information Technology,
King Faisal University, Hofuf, Saudi Arabia

*Abstract*—Symbolic AI is indispensable for the current LLM agents that are used for example to reason the context of the questions. An expert system is a symbolic AI that can explain the reasoning it reached to, which typically is a rule-based system has been attractive for different domains such as medicine, agriculture, and operations. On average, these systems involve hundreds of rules that are instable; moreover, they are coded at low levels of abstraction. Therefore, designing and reengineering an expert system is still costly and needs technical knowledge because of the manual process and maintaining of a low-level abstraction. On the other hand, model-driven architecture (MDA) has proven to be a successful technology that raised the abstraction level and formalized it to automate software development. It specifies business aspects in the platform-independent model (PIM) and implementation aspects in a platform-specific model (PSM). It then automates mapping between them using a standard mapping language called Query-View- Transform QVT. This paper argues that utilizing MDA principles such as the automation and abstractions represented by the descriptor PIM and PSM and mappings metamodels will not only overcome the instability of rules of expert systems, but also provides new insights for its usage. Therefore, this work proposes an MDA-compliant methodology that adopts a UML sequence diagram, a class diagram for the PIM descriptor, and a generic PSM) based on production rules. Moreover, a UML profile to support lacking features in the sequence model has been developed. However, the paper argues for a new kind of process-oriented expert system. Therefore, it not only allows domain experts to develop or participate in expert systems but also reduces the cost of developing new systems and re-engineering or maintenance of the critical and large-scale legacy expert systems.

*Keywords*—*Model-Driven-Architecture(MDA);Unified Modelling Language (UML); Platform-Independent Model (PIM); Platform-Specific Model (PSM); Query- View- Transform (QVT)*

## I. INTRODUCTION

Expert systems (ESs) are historically the most successful product of artificial intelligence (AI) [1]. It is widely used and designed to solve complicated problems that require reasoning about knowledge using mathematical logic. In fact, "based on Chatbot Agent—Google Bard" makes use of some reasoning based on logic to appear consistent and accurate. In AI, symbolic knowledge is represented as symbols that model concepts and relationships in the form of rules. It turns out that a variety of ESs that have been developed successfully and served stakeholders over a long period of time have become assets for many organizations. For example, an expert system was developed to provide clinical interpretations from thyroid hormone pathology tests decades ago. It had like 700 rules representing the knowledge-based approach, which provided 6,000 interpretations per year [2]. Also currently, the MD Anderson Cancer Center Expert System [3] helps oncologists make more informed treatment decisions. This system has a large knowledge base for oncology; it is expected to have thousands of rules.

However, the commonality among these systems and many others is the rapid change in the knowledge base because of the progress in the landscape of the field; for instance, new treatments, diagnostic techniques, and clinical guidelines that support the domain, as well as a technological change. Moreover, there is an essential business requirement for integrating these systems with others, such as electronic health records (HER), enterprise resource planning (ERP), and others, to increase their capabilities.

Although ES is an old field of study and there are many competing disciplines contributing to decision support problems, such as data mining, deep learning, and large language models (LLM), which are data-oriented approaches [4], ES conserves its unique properties of supporting problems that are rule-based and the capability of explaining reasoning. On the other hand, a hybrid model is typically used [5], whereby a weakness (i.e., learning capability from unstructured data) of one can be improved with the strength of the other (i.e., machine learning). In fact, rules are an intrinsic element of organizations, so decisions are driven by rules that support the production of services or products. Moreover, ES has attracted to some extent new domains, such as environmental data management analysis [6] and policy automation [7].

Given this situation and the fact that adopting code-based approaches such as Pyke, CLIPS, Prolog, Lisp, and other platforms for building ESs remains critical and costly, improvement is essential [8]. For instance, one reason for the interruption of some big systems (i.e., Garvan and IBM) is the high cost of maintenance and reengineering [9]. For example, frequent rule changes, driven by tech or business needs (i.e., adding some quality), require tedious hard coding, raising maintenance and re-engineering costs. On the other hand, there is a need to find new ways to make it easier for domain experts to develop or participate in ES. However, utilizing new engineering methodologies can provide a remedy for these problems.

However, the Object Management Group (OMG) has developed Model-Driven Architecture (MDA) as a methodology for automating software development, which is

standardized by OMG and supported by many tools, such as the Meta-Object Facility (MOF) design language [10], Object Constraint Language (OCL) [11], and XML Metadata Interchange (XMI) [11]. MDA encourages investment in metamodels using MOFs that are platform-independent, thereby paying back the low cost of development. A metamodel, or model of model, is a conceptual model of some design language that could have different implementations. In fact, having for different languages formally an equivalent representation (a phenomenon frequently referred to as syntactic sugar) is common (i.e., for language designers), such as Structured Query Language (SQL), relational algebra, tuple relational calculus, and query by example (QBE), which have the same underlying structure—the same metamodel [12]. In this case, this metamodel is known as abstract syntax, and its representation is called concrete syntax [13]. The strength of decoupling abstract syntax from concrete syntax allows much freedom in having different concrete syntaxes that re-use the same formal abstract syntax. Consequently, concrete syntax changes do not necessarily require abstract syntax to change. This flexibility allows you to re-use the same supporting software for a metamodel. For example, if we have the abstract syntax of SQL (a metamodel) itself and it is modeled in BNF (meta metamodel), it would be possible to send it a tool to render it into OCL, relational algebra, and probably first-order predictive calculus (i.e., specification of some constraints or assertions over a schema is needed). More importantly, MDA realizes the automation of mappings between these different metamodels using a standard mapping language called Query-View-Transform (QVT). Therefore, this separation of concerns, for instance, allows us to adopt a change without much cost incurred because of the high degree of sustainability. We argue that this trend is most suited to the instability of rule phenomena or the frequent changes in business requirements and technology in the context of ES.

### A. Expert System

An ES is a computer program that manipulates facts and rules that constitutes a knowledge of some domain to solve complicated reasoning problems efficiently and effectively [14]. These problems require domain experts'intervention to capture the knowledge [15] which is limited by human capability of handling hundreds of factors at the same time.

Many applications in health, industry, or education fields are using ES [16], [17], [18]. In these cases, the utilization of ES is geared toward delivering exemplary performance in addressing intricate challenges within a particular domain. ESs are instrumental in offering explanation and incorporating symbolic reasoning methodologies during problem-solving processes. Consequently, diagnostic ESs continue to hold prominence as the most frequently employed applications in this regard [19]. More important, ES can serve different class of problems such as acting as a *classifier, predictor, and estimator* to serve different sort of application domains that require automation support for decisions.

Moreover, ES has two core components: a knowledge base (KB) and reasoning engine [20]. The propositional ES has a KB formalized using propositions logic; it models the real world in the form of predicates and rules that can be evaluated to check its truthiness with initiations of variables. The KB is like a warehouse that contains knowledge about a specific domain captured by human experts in a form of production rules. Typically, it is a result of an expensive process called in literature a knowledge-Acquisition that involves strong communications between domain experts in one knowledge area and ES developers. It is the critical part in engineering ESs because of the requirement for developers to transfer this rigorous nature of rules of a domain knowledge into symbolic abstraction using logic-based structure. The classical engineering methods in this context follow an informal approach where engineers build informal models to capture the requirements of the system [21]. While the reasoning engine is an interpreter that draws a conclusion from premises that are represented using like first order predicate calculus (FOPC). The well-established theory behind that is the mathematical logic and theorem proving [22]. An example of the theories behind reasoning are Mode's pones [23] that allows to draw a conclusion from premises; It says if P proposition is known and has the fact that P implies Q then we also know Q as a conclusion. Traditionally may software tools that are aczt as inference engine are used to support the execution of ESs: Shells like Drools [24], G2 [25], JADE [26] which is based on the theory of agents and the standard FIPA [27], and Oracle Intelligent Decision Management (OIDM) [28]. These inference engines or Shells are classified based on forward changing (goes from known premises to reach goals by applying rules) or backward changing (works from goal to find the necessary premises) types of reasoning. Because users and domain experts need to understand how these tools reach to some conclusion, Shells have typically adopted a third essential component the explanation facility that could be in different forms: sequence of rules, conditions under which a rule fires and the conclusion it draws, and high-level detailed specification for the reasoning step [29]. In addition, the user interface component is for inserting quires, inputs and converting the rule from the internal representation to be user-understandable form.

### B. Model-driven Architecture

OMG has developed and standardized MDA with the aim of developing software without writing code. Models are first-class entities in MDA that enable the re-use of existing software assets, thereby reducing the complexity and cost of development. In MDA, a Platform-Independent Model (PIM) is used to specify the application concepts, while a Platform-Specific model (PSM) is used to specify the implementation issues independent of a technology. Then a standard mapping between PIM and PSM is specified using a Query- View-Transform (QVT), a standard mapping language that performs mappings between metamodels. MDA is an approach for building models, making the transformation of a source model into a target model [11]. To realize MDA, OMG has worked in the set of tools and standards that have been defined and standardized by OMG to support a good infrastructure. These OMG standards include UML, Meta- Object Facility (MOF), XML Metadata Interchange (XMI), Object-constraint language (OCL) and QVT. However, MDA has three types of abstraction: Computational Independent Models (CIM), PIM, and PSM. Each type is an abstraction technique for focusing on a particular part of concerns within a system and can be represented via one or more models.

A PIM is a conceptual model that is platform-independent and focuses on modeling domain concepts. PIM has a higher degree of independence from different platforms (e.g., .NET, CORBA and J2EE). In order to implement a PIM concept there should be a corresponding implementation abstraction involves a concept which can map it, typically the PSM abstraction. PSM is built from some technology's perspective but independent of it because MOF-based language is used. For example a developing database application requires to model the relational model using MOF; the example used in QVT standard document then by automating the transforming of the PIM instances into PSM instances, the main part of writing code is achieved. Thus, the PSM is considered a high-level APIs specification for a well-established platform such as database manager in. However, the main question here is how it would be able to relate a  PIM concept with PSM concepts that will be automated. Also, what are the suitable models for representing problem space and solution space?

*C. Query / View / Transformation*

The OMG has defined a standard for model transformations in the MDA architecture which is QVT. It represents the intrinsic activity in MDA engineering of applications whereby it converts a source model to a target model. It is required that the source model and the target model must both be compliant with the MOF meta-model [30]. QVT defines three specific languages named: (1) Relations, (2) Core (3) Operational/Mapping. Relations and Core are declarative languages with two different levels of abstraction. The QVT Operational / Mapping is an imperative language, it provides common constructions in imperative languages (i.e. loops, conditions.).

Because the relation language allows a round trip mappings between metamodels, this paper uses relational QVT language, specifically MediniQVT [31] engine to test the developed rules of the proposed mappings. Relational QVT has two main clauses: (1) Check only and (2) Enforce.

The "check only" clause focuses on validating source and target models against pre- defined rules and constraints without modifying the target model and act as precondition. On the other hand, the "enforce" clause carries out the actual transformations on the target model based on predefined relationships made between PIM and PSM metamodels. For instance the following example is a part of a complete transformation example that maps UML conceptual model (PIM) for database application into a relational model (PSM) [30] :

**Top relation PackageToSchema  {**

checkonly domain uml p: Package {name=pn}

 enforce domain rdbms s: Schema {name=pn}}

The provided example shows a transformation rule or relation named "PackageToSchema" that transforms objects from the "uml", the source domain to "rdbms" target domain. It says that if it is true that an instance of package exists in the source, creates a corresponding instance of type schema in the target. A domain is a typed variable that can be matched with a model of specific type which consists of patterns (i.e.   p:

package {name=cn}. A pattern can be grasped as a set of variables and constraints that needs to be bound by elements from a model to satisfy it with a valid binding. A domain pattern is a blueprint for the objects and their properties that must be found, changed, or made in a candidate model in order to meet the relationship [30].

## II. RELATED WORKS

The development of ESs using a model-based approach within the MDA process mapping is still an ongoing area of research, with limited studies focusing on this domain. Chungoorae et al. [14] suggested an approach based on MDA and ontology-driven system development to implement Interoperable Manufacturing Knowledge Systems (IMKS) for product lifecycle applications. It involves developing the PIM level using manufacturing core ontology and transforming it into a PSM in the XKS format. This approach lacks generalization for PSM and is limited to specific platforms. Moreover, it omits the mention of the mapping language used from PIM to PSM. Additionally, it primarily focuses on data-oriented ES type.

The BOM approach, within MDA, automates software generation using PRISMA architectural models, as described by Cabello et al. [15]. The approach utilized conceptual models, PIM, and CIM without specialized languages, resulting in generated program codes for C# .NET. This approach lacks generalization for PSM and is limited to specific platforms. Moreover, it omits the mention of the mapping language used from PIM to PSM. Furthermore, it does not specify the type of ES being utilized.

Yurin et al. [16] proposed using MDA to generate an ES for analyzing construction material damage. The implementation involves conceptual models, rule visual modeling language (RVML) for mapping PIM to PSM, and implemented in the form of a software prototype personal knowledge base designer (PKBD). This approach lacks generalization for PSM and is limited to specific platforms. Moreover, it utilizes an operational language that does not take maintenance into consideration. Additionally, it primarily focuses on data-oriented ES type.

Another research by Maylawati et al. [17] focused on UML diagrams (use case, class, and sequence) to describe ES components, including actor interaction and object relationships. Each use case requires a sequence diagram for normal and alternative processes, while the class diagram represents object interrelationships and adaptable attributes/methods for problem-solving. This approach involves the development of ES using UML diagrams in a general sense, without specifying the principles of MDA.

In study [18], the authors proposed developing decision-making modules for an intelligent system based on MDA principles. They start with the CIM, transforming spreadsheet data into a conceptual model using UML class diagrams. The PIM is then created as a rule-based module, formalizing decision tables with concepts from the CIM. RVML schemas represent these concepts. The PSM is developed depending on the knowledge representation language and converting decision tables into an RVML module. Program codes are generated

using PKBD for the specific platform. This approach lacks generalization for PSM and is limited to specific platforms. Moreover, it utilizes an operational language that does not take maintenance into consideration. Furthermore, it does not specify the type of ES being utilized.

Overall, the existing literature lacks comprehensive compliance with MDA principles. For instance, strong support for re-usability and platform independence such as having PSM that is special-case and PIM that is cluttered with PSM aspects. Also, the lack of maintenance considerations because of using operational mapping languages as well as in overall there is no support for it. Our proposed approach aims to generate a generalized PSM that can be adapted to different platforms as well as support maintenance. This is enabled by utilizing a relational QVT and developing reverse mapping rules. Moreover, this work recognizes a new type of process-oriented ES and enables both types of ES: process-oriented and data-oriented.

## III. RESEARCH METHODOLOGY

This paper aims to automate the creation of rule-based ES through the utilization of MDA. MDA, as explained, intends to automate software development without writing code by raising the abstraction level, so the process is driven by high-level models and mappings. The proposed approach argues for process-oriented ES as well as data-oriented ES. However, MDA principles compliance is a target for this work as well as considering the maintenance support. Therefore, these reforms the classic process of developing ESs that use a code-based approach into a new methodology for the development that uses models as first-class entities. The following subsections deal with the following questions: What are the suitable models to represent the PIM and PSM metamodels? Is there any gaps (i.e. concepts) exist in the source or target design language? How are PIM concepts related to PSM concepts? How to abstract Shells (platforms) in a generic PSM? And finally how low-cost maintenance is supported?

The answer to these questions can be organized around main principles of the proposed approach: (1) Modeling Expert Business Rules in PIM metamodel, (2) Building a UML Profile, (3) Developing PSM of Production Rules (4) Building the mapping between PIM and PSM.

### A. Model Expert Business Rules in PIM Metamodel

The expert system usually stemmed from business rules that represent the decision instrument which requires automation support. These rules classically captured manually in a form of conditional statements rendered as FOPL that probably augmented with the syntax of platform. They are two abstraction levels act as one level which makes it difficult for domain experts as well as developers to deal with these technical aspects. In contrast, the PIM is an alternative abstraction decuples the technical aspects of ES from application concepts therefore pertain to a conceptual model for problems under consideration typically developed independent of a platform so hides implementation details. Due to its abstract nature, it can capture essential features and requirements, which can help experts from different domains communicate by enabling shared understanding. MDA

approach proposes MOF-based language such as any UML models to act as PIM metamodel [10]. The PIM can be metalevel 1 or metalevel 2. The concrete instances of a UML model that represent specific case of ES (i.e. some diseases diagnose system), while UML is at metalevel 2 because MOF can model it which is a metalevel 3 [32]. This flexibility allows different modeling representation capabilities that address the diversity of system. The business rules for ES in code-based approach are dependent on a clear understanding of the ES's requirements. This understanding may be achieved through close collaboration with domain experts, stakeholders, and end-users. But usually there is no corresponding explicit formal model such PIM metamodel that can capture these requirements.

More importantly this work argues for process-oriented that centers on the workflow or operational procedures of a given expert system, such as production processes, scheduling operations, and generally processes. In this type of systems, the rules are injected within a process, not like the classic ES that focuses on data so called data-oriented, process is dominating element. For example, in an academic system; a student can only register for 12 credit hours if his GPA is less than 2 and, while in manufacturing system; IF machine temperature reaches critical threshold, THEN stop the machine and call maintenance procedure. The former rule is a constraint augmented with an action – register, while the later expresses object properties constraints, the temperature of the machine. The distinction between them is useful for acquisition of knowledge process as well as it has impact in the design under the context of this work. Although, the model we have introduced is capable of accommodating either of these types, its emphasis is primarily on the latter (i.e., process-oriented). However, business rules, are set of guidelines or policies that dictate how an organization operates, are frequently embedded within processes to ensure and enforce consistency and integrity [33]. It is therefore a good candidate for PIM.

However, this paper proposes specifying business rules for ES in PIM, using both class and sequence diagrams. In which class diagram serves as a descriptor of facts or data with their schemas that is needed by the ES, while the sequence diagram is used for capturing the behavior of the system that usually is augmented with rules. Although there are alternatives to this design decision such as activity diagram, sequence diagram is less elaborative so more compact and easier to learn and communicate the problem of ES. Nevertheless, the relationship between UML sequence diagram and UML class diagram is that the sequence diagram involves objects and messages that are eventually comes from classes that are supposed to be fully specified by the UML class model. Therefore, the UML class diagram supplies the sequence diagram with schemes of data and their relationships.

*1) PIM sequence diagram metamodel:* A metamodel is a model of the model that is required here to capture the elements needed to support process-oriented ES instances; that is a UML PIM metamodel. We need to investigate the big enough and suitable UML sequence diagram metamodel to be ready for representing the instances of PIM that will be developed by domain experts and developers as well as guides

the mapping process later. A developer in this case typically develops models for ES at metalevel 1 where it models objects of specific ES such as a manufacturer production expert system or car faults diagnoses system. The intention of a sequence diagram in UML is to communicate the specific behavior by sketching the sequence of messages communicated between objects in a system so it's a dynamic view. Fig. 1 presents the necessary and big enough metamodel elements needed to model any business rules for experts.



Fig. 1. PIM sequence diagram metamodel [34].

Fig. 1 is the abstract syntax of the sequence diagram that is part of the UML metamodel which fits the problem addressed by this work. The metamodel says a Message (with name and kind properties) can have zero or more Arguments (with name, direction, value properties). It can have a Return Type that specifies the type of value that is returned. In addition, the Message has the Message End indicates two ends of a message exchanged between two Lifelines (usually named element) where the first represents the source (base) and the second is the destination. A Message End is either an Object or an Actor. A lifeline can have more than one message. It indicates that a message has been successfully transmitted from the sender to the receiver, and that the receiver has finished processing the message.

The CombinedFragment is an essential component in diagrams that allows for the specification of complex control structures such as loops, parallelism, and conditions. Its behavior is determined by the chosen InteractionOperator, which dictates how the fragment behaves. For example, when using "alt" as the InteractionOperator, the CombinedFragment represents a choice of behavior where only one option is executed. Semantically an Operand selection within the "alt" fragment is based on guard expressions, which determine the conditions for executing each option. The use of the "else" guard expression represents a negation of all other guards within the CombinedFragment. If none of the options have guards that evaluate to true, none of them are executed, and the remaining part of the diagram continues. Note that an enumeration class called InterationOperKind is used to supply the kinds required for InteractionOperator. It is a useful feature can be utilized in expert systems in exceptions as well as normal conditional state.

*2) PIM class diagram metamodel:* The UML class model utilizes class diagram notation [35] to describe the data and their relationships using class, properties, inheritance

(generalization-specialization) and association concepts. For example, in the journal system a Reviewer and Paper are classes (objects of metalevel 1), and a Review class has the association with Paper between them that represents the relationship a reviewer provides a review for a certain paper which also has an association with a Review class that represents the feedback. In the following, Fig. 2 illustrates the UML class diagram metamodel that is part from UML standard specifications developed by OMG to act as a descriptor for data and facts needed in ES with their integrity constraints.



Fig. 2. PIM class diagram metamodel [30].

Fig. 2 shows part of the abstract syntax of the UML class diagram [20]. In which UMLModelElements are classes, interfaces, packages, and relationships. For example, a class diagram could include several UML model elements such as Classes, Attributes, and Associations, to model classes of objects with properties and the relationships between different objects in a system. A metaclass class in Fig. 2 is a central concept of the language and it is a kind of a classifier. A classifier in UML means a class that can be instantiated or has instances different from Abstract classes that do not have. Classes are defined by a set of attributes and methods that objects of that class will have. Because usually a class like Reviewer has attributes: ID, name, and Area of interest there is of metaclass Attribute to model this. Attributes are characteristics or properties of a class, and they define the data that objects of that class will have. In addition, Association, which describes the relationships between two or more classes. Associations define how objects of one class are related to objects of another class. Moreover, a Package is a grouping mechanism used to organize related elements, including classes, interfaces, and other packages. A PackageElement is an abstract class that is a kind of package.Therefore, a diagram usually exists in a package. Another critical UML element is the Classifier, which describes a set of objects that share common characteristics and behavior. Lastly, UML includes PrimitiveDataTypes, which are basic data types built into the modeling language or programming language used to represent

data. Examples of UML primitive data types are Boolean, integer, and string.

*3) UML profile for PIM metamodel:* A profile is a powerful lightweight extension mechanism that adds concept, syntax, and semantics to the metamodel [13]. It does not require substantial change to the metamodel so less costed approach because UML editors does not change because of extensions. A profile consists of stereotypes, tags, and metaclass classes of the elements that need to be extended, which are classifiers. The stereotype represents a new concept or syntax that is needed for extension while a tag adds some properties if needed to the stereotype. However, it is necessary to have an end of an extension that represents some metaclass class; a solid line notation designates this extension usually drawn between the two ends in UML standard. Profiling plays a pivotal role in filling the gaps in the metamodel. Because in this work there is a gap in the Sequence Diagram; it is not defined by UML Sequence metamodel which are the modeling elements: Not, OR, and Head that exist in the target (Shells). Therefore, there is to design a profile in order to allow developer/domain expert to use these concepts. Fig. 3 shows the UML profile diagram for the proposed approach.



Fig. 3.    PIM profile diagram metamodel.

In Fig. 3 a profile is designed for the PIM sequence diagram with three distinct stereotypes needed. These stereotypes serve to provide concepts that the PIM sequence model is lacking. It basically extends two metaclass classes: Message and Actor because semantically OR and NOT are related to messages in the level of abstraction and later to predicate as we will see in the mapping. The first stereotype of the 'Actor' element is named a 'head,' concept to allow designers of ES to utilize the concept of specifying the focal point which gets the benefit of the service provided by sequence diagram. A tool cannot easily determine without cost, so the profiling mechanism is a less costed solution.   The second stereotype pertains to messages and introduces an 'OR' operator concept. This 'OR' is an indicative of multiple possible interactions that can exist between messages, a common practice in ESs (multiple rules with same head). The

third stereotype extends the Message to introduce a 'NOT' operator concept. In this context, the 'NOT' signifies conditions or interactions that are explicitly negated or excluded (it is also common in ESs). Therefore, the utilization of these stereotypes within the profile diagram serves the purpose of specifying rules that involve disjunction, negation and as well discriminating the Head of the rule. We are ready now to look at how we design a generic PSM.

### B. Develop PSM of Production Rules

A production rule traditionally used in different fields such complier, natural processing languages and logic which specifies how input stimuli are transformed into output responses (produce a symbol output from a symbol input). It consists of a set of rules, each consisting of a condition and an action or LHS and RHS. A condition specifies a set of constraints that must be satisfied by the input, whereas an action specifies a set of operations to be performed on the input [36]. To fire it means to replace the LHS with the RHS. To model production rules using PSM, it is necessary to identify the specific elements of FOPC [37] because it is what Shells of ES are based. Fig. 4 illustrates the PSM metamodel for the proposed ES.



Fig. 4.    Rule-based expert system PSM metamodel.

In this context, a typical production rule might be expressed in natural language, as in the example provided: "If the car won't start and there is no clicking sound when the key is turned, then the problem is likely a dead battery". To get the metamodel of FOPC as in Fig. 4, since after investigation instances of rules exist in this world, we need a metaclass called Rule in the PSM metamodel. A Rule can be identified by ID, so we need an attribute called ID of type integer. Since a Rule has ascendant that act as a set of conditions must be met for a rule to be applied or executed, we call it the right-hand side and the consequent or action    will call it left-hand side. Therefore, it is necessary to define a metaclass classes called RightHandSide(RHS) and LeftHandSide (LHS) because there are many instances will be of this kinds for rule. Moreover, we observe that each of them consists of a basic building block known as predicates; so, we need to model a metaclass class called Predicate because there are many instances of it in a single rule. In addition, usually a Predicate involves parameters that are variables that should be bound during execution; we need a model it as a metaclass class called Parameter, which provides parameter's name and type. Now we can turn into the relationships between the PIM and PSM.

## C. Build the Mappings between PIM and PSM

The aim of MDA eventually is to map the source (PIM metamodel) into a target (PSM metamodel) which acts as a part of writing the code in the development process. This model transformations should be done through the standard language QVT which is independent of both ends. In the context of rule-based ES, the expected output of this mapping process is a PSM instances that can be used to represent executable facts and rules for the target platform. It is typically a UML model instances or objects of metalevel 1 [32]. The PSM should be able to implement the PIM objects so translate them into specific constructs or patterns that are suitable for the target platform. The following subsections discuss the mapping rules of the proposed approach: (1) PIM Sequence Metamodel to PSM Metamodel, (2) PIM Class Metamodel to PSM Metamodel.

*1) PIM sequence metamodel to PSM metamodel:* Sequence diagram is utilized initially to facilitate the decision-making process by depicting the business process aspects of the ES in a high-level model that allows domain experts to communicate the problem easily. However, the goal is to convert business rules that act as the knowledge of expert in some domain into implementation using PSM concepts. As consequence of this, we need to find relationships between sequence diagram concepts and the PSM concepts which are production rules that are commonly represented using FOPC [37]. Table I shows these relationships of both metamodel concepts.

TABLE I.    PIM Sequence Diagram to PSM Rule-based ES Mapping Rules

| Rule | Transformation Rule | Source Model | Target Model |
|------|--------------------|--------------|--------------|
| R1 | MessageToLHS | Message: kind= 'goal' | LeftHandSide |
| R2 | MessageToPredicate | Message: kind= 'normal' | RightHandSide |
| R3 | AltToHead | Combinedfregment: interactionOperator='alt' and Message:kind= 'goal' | LeftHandSide |
| R4 | AltBodyToRHS | Combinedfregment: interactionOperator='alt' and Message:kind= 'normal' | RightHandSide |
| R5 | ArgumentToParameter | Argument | Parameter |
| R6 | MessageEndToRelationship | MessageEnd | Relationship |
| R7 | Negation | Message (has NOT) | Negated predicate |

These transformation rules play a crucial role in converting from PIM sequence model to PSM of rule – based ES. For instance, the "MessageToLHS" rule maps a message with the 'goal' kind into the LeftHandSide format, while the "MessageToPredicate" maps a message with the 'normal' kind into the RightHandSide. The "AltToHead" rule transforms a Combinedfregment with interactionOperator='alt' and message with kind= 'goal' into the LeftHandSide. The "AltBodyToRHS" rule transforms a Combinedfregment with interactionOperator='alt' and message with kind= 'normal' into the RightHandSide. Similarly, the "ArgumentToParameter"

rule converts an argument into a parameter, and the "MessageEndToRelationship" rule transforms a message end into a relationship. Lastly, the "Negation" rule is employed when a message includes 'NOT' to create a negated predicate in the target model. These rules are acting as separate artefacts so preserve the separation of concerns principle.

To automate the mapping process, mapping rules must be specified using QVT standard transformation language. The QVT mapping rules has the following structure:

transformation map (source: sequence, target: psm)

This transformation specification is like a procedure map or make a transformation from a source model represents as a sequence diagram (source) to a target model represents a PSM (target), it is telling the tool that the mapping direction.

The following formalizes the informal mapping rules specified in Table I for the specific mapping between PIM and PSM. Rules will be numbered (ascending order) for easier reference (All these rules are tested using (MediniQVT tool).

R1:

top relation MessageToLHS {

k, cn: String;

checkonly domain source m: sequence::message{kind = 'goal', name =k};

enforce domain target p:psm::predicate{at= a:psm::LHS{}, name = k};

where {ArguementToparmeter(m,p);

In R1, the source domain, involves a pattern that checks for messages in the sequence diagram with a specific *kind* = "goal" and a *name* that will bound using variable "k" based on instances of the source. If it's true, then in the target domain, the rule enforces the creation of a *predicate* in the PSM; with a repository, with the name as "k" and an attribute "at" assigned to a parameter "a" of type "psm::LHS". Additionally, the rule includes a constraint that the relation "ArgumentToParameter" must be called after executing this rule which ensures transformation from argument to parameter. Because mapping to arguments is postcondition (executes after the first part above) and a complement, *where clause is added.*

R2:

relation MessageToPredicate {

pn: String;

checkonly domain source m: sequence::message{kind='normal', name=pn};

enforce domain target p : psm::predicate {at = a :psm::RHS{}, name=pn};

where { ArguementToparmeter(m,p);}}

The rule R2 directs a transformation from the sequence domain to the psm domain. It checks if there is a *message* with a specific *kind* (normal) and *name* = pn (bound to the current instances in the source) in the sequence diagram. if true, then in the target domain, it creates a *predicate* with a matching

name(pn value) and assigns to the attribute "*at*" a parameter "a" of type "psm::RHS". The rule also includes a constraint, postcondition that the rule R5 -"ArgumentToParameter" for parameter transformation must be executed afterwards.

R3:

top relation AltToHead{

cn,n:String; o: Integer;

checkonly domain source f:sequance::combinedfregment{

interactionOperator='alt',ID=o,lifelinee=k:sequance::LifeLine{name=n,

messages=m:sequance::message{name=cn,kind='goal'}}};

enforce domain target lt:psm::RHS{ID=o};

enforce domain target ltt:psm::LHS{ID=n};

enforce domain target p:psm::predicate{name=cn,at=a:psm::LHS{iD=n}};}

R3 is a rule that checks if the sequence model has a fragment with an interaction operator set to "alt" and an *ID* matching "o" (bound with current instances).Also, it verifies a *lifeline* with the *name* "n" and a *message* with the *name* "cn" and kind "goal" within that fragment. If true, then in the PSM model, the rule enforces creation for the right-hand side (RHS) with an *ID* matching "o", a left-hand side (LHS) with an *ID* matching "n", and a predicate with the *name* "cn" and *an* attribute referring to the left-hand side with ID "n." This rule does the initiation task where the rest of the rule will base.

R4:

top relation AltBodyToRHS{

i:Integer;cn,n:String;

checkonly domain source f:sequance::combinedfregment{

{iD=i,interactionOperator='alt',lifelinee=k:sequance::LifLine{name=n,

mesages=m:sequance::message{name=cn,kind='normal'}};

enforce domain target p:psm::predicate{name=cn,att=a:psm::RHS{iD=i}};

where { ArguementToparmeter (m,p);}}

R4 is a complement rule to R3 rule, asserts a *combined fragment* in the source sequence model with an interaction operator set to "alt" and an ID matching the given value of "i". It also asserts a lifeline with the *name* "n" and a *message* with the *name* "cn" and the *kind* "normal". If true, in the target domain, creates a *predicate* with the name "cn" and the attribute "att" assigned the right-hand side (RHS) object with the ID value for "i".Also, R5 is postcondition so needs to be executed afterwards as where clause exists.

R5:

relation ArguementToparmeter {

Cn ,n,q: String;

checkonly domain source m: sequence::message{kind = q,name = Cn,

pars = w:sequence::argument{name = n}};

enforce domain target p:psm::predicate{name = Cn,

args = k:psm::parmeter{name = n}};}

R5 relation asserts for a message in a sequence diagram with a specific *kind* and *name*, and pars attribute with *argument* that has *name= Cn, if turr in* a target domain, will be the creation of a *predicate* with the same name and parameter.

R6:

top relation MessageEndToRelationship{

checkonly domain source b:sequence::messageend{name = cn,

sender = e:sequence::message{ kind = 'normal' }};

enforce domain target w :psm::relationship{name ='AND',

srcP = ps:psm::predicate{}};

where {MessageToPredicate(e,ps);}}

R6 asserts if a *message end* in the source sequence diagram with a specific *name* and a *sender* that is a "normal" message. Accordingly, in the target domain, it enforces the creation of a relationship with the *name* 'AND' and a *source* predicate. This rule builds the relationship between predicates that usually is 'And' if not specified 'OR'. R2 is required as postcondition.

R7:

top relation negatedTopredicate {

Cn ,n: String;

checkonly domain source a:sequence::message{kind= 'negated',name = Cn};

enforce domain target o:psm::predicate{name=Cn +'not'};}

R7 rule checks for a message in the source sequence diagram with a specific kind "negated" and a name matching the variable "Cn". If true,in the target domain, it enforces the creation of a *predicate* with a name formed by appending "not" to the original name.

*2) PIM Class Metamodel to PSM Metamodel:* The class diagram is utilized initially to act as a descriptor for data and facts needed in ES. However, the goal is to convert the facts of experts in some domains into implementation using PSM concepts. As a consequence of this, we need to find relationships between class diagram concepts and the PSM concepts. Table I shows these relationships after a close investigation of both metamodel concepts. Table II shows transformation rules representing the mapping between PIM class diagram to PSM rule- based ES.

In Table II, the "ClassToFact" rule performs the transformation of classes into a fact in the target model. Similarly, the "AttributeToParameter" rule is employed to convert an attribute into a parameter. These rules play a vital role in the process of adapting and reshaping data within the

modeling contextualizing the need for facts that represent the essential part of knowledge.

| Rule | Transformation Rule | Source Model | Target Model |
|------|---------------------|--------------|--------------|
| R1 | ClassToFact | Class | Fact |
| R2 | AttributeToParameter | Attribute | Parameter |

The QVT mapping rules starting with this statement:

transformation map (source: class, target:psm)

This transformation specifies the mapping between a source model represented as a class diagram (source) and a target model represented as a PSM (target).

The formal QVT mapping rules corresponds to Table II:

R1:

top relation ClasstoFact{

 Cn , n : String;

checkonly domain source a:cla::classs{name = Cn};

enforce  domain target o:psm::Fact{name = Cn};

where {AttributeToParameter(a,o);}}

R1 relation transforms a source model (class) to a target model(psm). It checks for a *class* in the source domain with a specific name "Cn". In the target domain, it enforces the creation of a *Fact* with the same name "Cn". Additionally, it includes a constraint R5 executes afterwards that ensures the mapping of attributes from the source class to *parameters* in the target for *Fact*.

R2:

relation AttributeToParameter{

Cn , n , v: String;

checkonly domain source a:cla::classs{name = Cn , attribute = ar:cla::Attribute{

name = n, value= v}};

enforce  domain target o:psm::Fact{name = Cn, parmeters = w : psm::parmeter {

name = n, value= v}};

R2 relation is part of the "ClasstoFact" transformation. It checks for an attribute within the source class that matches the variable "n" and has a value matching the variable "v". In the target domain, it enforces the creation of a parameter within the Fact with the same name and value. This relation ensures the mapping of attributes to parameters during the transformation process.

*D. Round- Trip Mapping*

The extant issues encountered in the development of ES via a code-based approach has strong resolution through the contemporary application of MDA. This approach enables developers to focus on the high-level concepts of system

design, reducing the complexity of development and facilitating the reuse of code and knowledge [12]. For instance, MDA offers cost-effective maintenance through the utilization of automation and the implementation of a round-trip mapping mechanism. This work argues for support of maintenance using relational QVT language (QVT-r is supported by EMF). It reflects the changes that happened to PSM such as having a new version of the software of shells. More important the changes in the PIM model (i.e., business rules change) will not change the PSM or mapping assets so can be re-used. This adds great value to the re-engineering effort required for ES. By structuring mappings in this manner, developers gain enhanced and ease mechanism to make change such as update to the rules. This approach optimizes the maintainability of the system, as it streamlines the process of rule manipulation and adaptation within the MDA framework. More importantly, there are legacy ESs serving organizations for a long time that can benefit from this model whereas in extreme cases models can be reverse engineered so can be changed and synchronized automatically with required changes. For instance, a rule can be developed for round- trip mapping for PSM – PIM Sequence Diagram where each represents different sort of changes:

RM1:

relation LHSToMessage{

k : String;

checkonly domain p :psm::predicate{ at = a:psm::LHS{}, name = k   };

enforce domain target m: sequence::message{ kind = 'goal', name = k};

where {parmeterTOArguement (p,m);}}

The purpose of this rule is to enable the redirection of mapping from the left-hand side (LHS) back to the message, facilitating any necessary changes.

RM2:

relation PredicateToMessage{

pn: String;

checkonly domain target p: psm::predicate {at = a :psm::RHS{}, name = pn};

enforce domain target m: sequence:message {kind='normal',name = pn };

where { parmeterTOArguement (p,m);;}}

The purpose of this rule is to enable the redirection of mapping from the predicate back to the message. These rules establish a mechanism for reverse mapping that is not only applicable to the specific rules but also to other rules derived from the PIM Sequence Diagram to the PSM.

The round- trip mapping for PSM- PIM Class Diagram:

RM3:

top relation FactToClass{

Cn, n : String;

checkonly domain source o: psm::Fact{name = Cn};

enforce domain target a: cla::classs{name = Cn}

where {ParameterToParameter(o,a);}}

The purpose of this rule is to enable the redirection of mapping from fact back to the class.

RM4:

relation ParameterToAttribute{

Cn , n , v: String;

checkonly domain source o: psm::Fact{name = Cn, parmeters = w: psm::parmeter {name = n, value= v}};

enforce domain target a: cla::classs{name = Cn , attribute = ar:cla::Attribute{name = n, value= v}};}}

The purpose of this rule is to enable the redirection of mapping from parameter back to the attribute. Also, these rules establish a mechanism for reverse mapping that is not only applicable to the specific rules but also to other rules derived from the PIM Class Diagram to the PSM.

## IV. REENGINEERING AND NEW INSIGHTS

This section analyzes and discusses the feasibility, insights, and opportunities of using MDA at different scales of change in the scope of the ES. The reengineering of a system typically involves a radical change for the entire system to achieve some result, while maintenance tackles parts of a system to improve it by making corrective action or minor modifications [39]. In this work, we use the term reengineering in a broad context. For instance, maintenance could be applicable to any part of PIM, PSM, or mapping rules, while the process of making radical changes (from a code-based approach to a model-based approach) to the legacy ESs by using MDA is a re-engineering process. However, the capability of interoperability is also one of the main concerns of MDA, which is defined as the ability to seamlessly integrate different systems or components to exchange information and work together [11].

There are many reasons why current ESs need to change under the umbrella of maintenance or reengineering, for example, the need to interoperate or integrate with other systems. The basic assumption of the data or facts underlying ESs is to be provided in a static way for reasoning. Nowadays, this is not the case where data should be updated by dynamic systems such as in the medical field by EHR (i.e., supply patient data) or general business ERP (i.e., provide production information such as a master or detailed schedule). The data involved in such systems is not only current but also comprehensive. For instance, patients with new symptoms or a production machine show new odd behavior in one manufacture. Based on the application requirements, data needs to be pulled or pushed from these systems to the relative ES. It is obvious that manual pulling or pushing is not practical in this sense. This visibility is a sort of strong business requirement that must be achieved today. On the other hand, rules as shown are subject to change due to different reasons, such as progress or a shift in the landscape of the knowledge of a field (i.e., medicine), but we argue that our approach enables automated supplementation of rules because the transformation

process in MDA is a separate and dynamic process with stable mapping rules. This will not only provide a dynamic way of running the ES but also provide new insights. One can imagine that GPT agents (like ChatGPT or Bard) or similar intelligent systems can utilize this feature. Indeed, these GPT-based AI tools use symbolic knowledge-based questions or queries to find the relevant information or identify the context of the question (using inferencing rules and knowledge). This view suggests that an integration mechanism allows data to be outsourced as well as rules, so ES can be provided as a dynamic service.

In addition, interfaces of an old legacy systems became obsolete and so there is tendency to be upgraded to new standards such as Jetpack Compose [40] developed by google for building native Android applications, SwiftUI [41] for IOs and mac applications, CSS frameworks than enable web-access for ESs, and many others taking into account in single system such as mobile you find different of GUI standards.

We now turn to the question of how MDA can support this strong demand for integration with these different systems. On the one hand, MDA has the abstraction of a PSM that is based on MOF to represent the technical aspects of a platform, one of which is GUI platforms or others such as APIs for specific platforms (i.e., EHR and ERP). So specifically, PSM for any of this need to be developed and a couple of transformations using like QVT [38]. More than a decade ago, typically the integration between systems followed standards such as service-based system technologies, web services, and WSDL [42], JOSON and RESTFul [43] or SOAP [44]. They contributed to interoperability between different tools and encourage more integration to be practiced. Cloud systems are basically complex, diverse systems that use these standards. More importantly, the literature is rich with some of these standards that exist as PSMs and can be re-used to reduce development costs. On the other hand, PIM is a business-level abstraction built independently of even PSM, so the portion of PIM related to interoperability, such as GUI or others discussed above, can be projected using the transformation capability of MDA. In this case, mapping rules only need to be changed if the PSM is already published (GUI, WSDL). However, there might be intermediate steps (pre-processing steps) needed, such as using the QVT view maintenance [30] capability to map PIM into more refined PIM or PSM into more refined PSM.

To conclude, MDA has rich architecture support for change, such as interoperability, that can allow even non-MDA systems to integrate in a manner that reduces the re-engineering cost. More importantly, this interoperability in this context provides new insights into using traditional expert systems, such as exposing expert systems as a service, and gains the power to of dealing with dynamic changes in facts or the instability of rules.

## V. EVALUATION USING CASE STUDY

An academic advising ES has been introduced to bridge the gap between students and advisors by shifting advising, complaining, evaluating, and suggestions from traditional ways to a more contemporary one [45]. The decisions need to be made by students during their academic journey such as course

enrolment, course withdrawal, postponing study, etc. In this paper, we take on a scenario of a rule-based ES for academic advising in the university system. The need for ES for academic advising is to take a decision for different actions involves uncertainty and a couple of factors need to be tested. For example, the decision to drop a course for low GPA students has different consequences which is not straightforward decision. Similar thing can be said for postponed study, drop a semester and so on.

The ES will build the work plan by identifying the student through some important points:

- Perquisite courses.
- Knowing the student's performance.
- The student's weakness points.
- Domain skills.
- Skills that a student needs to improve.
- Student goals.
- Track the student pathway.

The ES works to facilitate the communication between students and the advisors by raising the student's performance giving some recommendations that help the low GPA student to develop specific skills for different semester actions such as course enrolment, course withdrawal, postpone study, etc.

*A. Developing Process Model for ES*

A student who is struggling with a low GPA might approach their academic advisor for assistance in considering the option of dropping a course for the current semester. The ES is responsible for determining when the low GPA student is eligible for dropping a course according to the following conditions:

- **Perquisites Course:** Perquisites course must be 'Not Major' category.
- **Skill Assessment:** Students skill must be a 'Weak Skill' in this course.

The advisory academic rules are:

- **Rule:DropCourse(SID,CID)=**IfGet_PreCourse(CList) AN Check_PreCouese_Category(CID)AND Check_Skill(Skill)
- In all other cases, the system does not allow the student to drop the course.

These rules consider the relevance of the PreCourse category and the strength of the student's skill set. If the conditions are met, the system facilitates the selection of appropriate PreCourse categories and skills, while disallowing the student from dropping the selected course. Conversely, if the conditions are not met, the system permits the student to drop the course if desired. Fig. 5 illustrates the sequence diagram outlining the process for dropping a course for a student with a low GPA.



Fig. 5. Drop a course UML sequence diagram.

Fig. 5 presents the process involving the Advisor, System, and PreCourse, Category, and SkillSet objects, illustrating how t odecide on the drop of a course. The interaction commences as the advisor engages with the system, focusing on the rules associated with dropping a course. The first message, "If the PreCourse category is classified as 'Not Major' or the student possesses a 'Week' skill in a particular subject (skill='Weak'), then the system proceeds with the selection of the PreCourse category (CID) through the predicate SelectPreCourseCategory(CID) and the selection of the skill (CID) through the SelectSkill(CID).Following this, the system responds with the message "DropCourse" for dropping the course, provided the conditions specified in the previous message are met. However, in all other cases, the system responds with the message "NotAllowtoDropCourse".

As noted earlier, MDA develops rule-based ES by mapping transformation from PIM to PSM. The mapping process is performed using the MediniQVT tool, where the source file is the Ecore file [46] representing the Sequence diagram PIM metamodel, and the target file is the Ecore file representing the Production rule PSM metamodel. Also, an XMI file acts as an input containing instances of the sequence diagram metamodel for this case for example is utilized. Subsequently, the relational QVT mapping rules mentioned above are applied to create the production rules of ES. Table III shows the mapping rules used and Fig. 6 shows a sample of execution for final result of mappings.

*B. Developing Data Model for ES*

As mentioned, the UML class diagram is used as a descriptor to represent the data and facts of ES. Fig. 7 explains the UML class model by the using example of a case study of academic advising system in university.

As shown in the UML class diagram is that Student has a relationship with Course. In addition, the Course has a specific domain (such as Math, programming) consisting of skills needed as outcomes for the course(s) in this domain. Further, a course sometimes has Prerequisite course; the association between Course and Prerequisite course, which models this business rule. This will enable checking the integrity constraint that that Student must take the Prerequisite course and pass it

before registering in a new course. According to the types of courses, there are two types: 1- Taken Course, 2- Next Plan Course. Taken course refers to the taken courses in the semester, and next plan courses refers to the planned courses in the next semester. Nevertheless, Student must have a study plan to follow according to the program requirements. The academic advisor wishes to help students complete this study plan successfully with low risk by making the right decision at the right time which is the source of the calling the experience of the ES.



Fig. 6. PIM sequence diagram for drop a course target file.

Fig. 7. PIM Sequence Diagram to PSM Rule -based Expert System Mapping Results

| No. | Predicate | Parameter | Predicate type | Justification |
|-----|-----------|-----------|----------------|---------------|
| 1 | Drop Course | Student ID, Course ID | Left Hand Side | This predicate represents the LHS predicate of the drop a course rule. |
| 2 | Get Pre-Course | Course List | Right Hand Side | This predicate represents the RHS predicate of the drop a course rule. It has an AND relationship with next RHS predicate |
| 3 | Select Pre-Course Category | Course ID | Right Hand Side | This predicate represents the RHS predicate of the drop a course rule. It has an AND relationship with next RHS predicate |
| 4 | Select Course | Course ID | Right Hand Side | This predicate represents the RHS predicate of the drop a course rule. |

As noted earlier, MDA develops rule-based ES by mapping transformation from PIM to PSM. Where the source file is the Ecore file representing the (Class diagram metamodel), and the target file is the Ecore file representing the (PSM metamodel). Additionally, an XMI file containing instances of the class diagram metamodel is utilized. Subsequently, the relational QVT mapping rules mentioned above are applied to create the facts of ES. Table IV shows the results of the mapping process, there are eight facts generated, each of which is associated with a specific parameter. These results provide a comprehensive description of the facts utilized by the domain experts in leveraging the ES effectively and Fig. 8 shows a sample of execution for final result of mappings.



Fig. 8. Academic advising system UML class diagram.

TABLE III. PIM CLASS DIAGRAM TO PSM RULE TO BASED EXPERT SYSTEM MAPPING RESULTS

| No. | Fact | Parameter |
|-----|------|-----------|
| 1 | Performance | Full Load |
| 2 | Student | Student ID |
| 3 | Study Plan | Credit Hours |
| 4 | Course | Course Name |
| 5 | Perquisite Courses | Skills |
| 6 | Low GPA Student | Current GPA |
| 7 | Taken Course | Date |
| 8 | Next Plan Course | NA |



Fig. 9. PIM class diagram for advising low GPA students in university target file.

## VI. RESULTS AND DISCUSSION

The investment on the quality of software development became evident that will pay back the cost. Having the case that many legacy ESs contributing to different domains exist over a long period, such as in medicine ,health, and education [16],[17],[18], necessarily entails requirement changes such as in platform or business rules. More importantly, the discovery of the entrance of ESs into new domains (environmental management and cybersecurity) requires flexibility and less costly development methods. However, using MDA in this work provides these qualities. The PSM metamodel from production rules developed as a generic PSM and mapping rules can act as assets so they can be re-used with the development process of any kind of ES; therefore, the principles of reusability, platform independence, are achieved and hence reducing the cost. For instance, changing Prolog with the Pyke platform for any reason such as utilizing a forward chaining tool instead of backward changing tool, does not cause a change in the PIM, PSM, and properly minor change to mapping rules. Also, changing to a new version of a platform such as upgrading to acquire new features, the proposed approach does not require to change PIM or PSM and mappings.

Similarly, in more extreme maintenance cases where rules are updated or modified, only the PIM (model instances) needs an update, the rest will be re-used therefore coping with the rules instability. Moreover, the raising of abstraction afforded by MDA, such as in the PIM descriptor, allows domain experts to participate or write expert system, which bridges the gap between domain experts and developers.

On the other hand, the integration of an ES with other systems (i.e. HER) under reengineering process or maintenance, is an inexpensive approach because of the re-using utility provided by metamodeling and formalizations (using MOF) of the descriptors: PIM, PSM, and mapping rules. Thereby provides costless reengineering. Because different Shells have different features, the PSM developed is standard one and therefore comply with the principle of platform-independence so like portability can be achieved. It can be modified to incorporate additional features if is to put into practice, but it should be the commonality among all shell platforms. The XMI standard allows either PIM or PSM to be migrated to another tool so can be edited or manipulated.

The sequence diagram, in reality, reflects the nature of interactions involving the business rules of a desired ES. However, this study argues for a type of ES that is process-oriented, where a set of actions with a sequence that represents constraints such as pre-conditions and post-conditions need to be specified for the desired outcome. For example, the process of checking ripe and unripe fruits, the process of optimization such as in manufacturing (i.e., efficiency of steal production), control process, real-time recommendations process, and planning and scheduling processes.

## VII. CONCLUSIONS

This work is about automating ESs from high-level models using the principles of MDA. ES is a long-sounding successful product of AI but lacks advanced methods of development and re-engineering, which leads to an increase in the cost of maintenance and development. Moreover, effective communication between developers and domain experts is a crucial yet challenging aspect of designing ESs. The inherent differences in technical knowledge and domain expertise often lead to communication gaps, hindering the accurate translation of domain knowledge into functional system components. However, MDA raises the abstraction level of the development of an application as well as provides a structured approach for automation so ES applications can leverage this feature. MDA decouples application concepts or domain of the problem that needs to be specified in the PIM metamodel from the technical aspects of implementation, which will be specified in the PSM metamodel; then mapping the first end (PIM) into the second (PSM) using the standard mapping language, QVT.

The proposed approach addresses some limitations in the literature, such as the lack of generic PSM and specific compliance to the MDA principles, as well as recognizes and supports a class of expert systems identified as process-oriented ES. A UML sequence diagram is used to model business aspects of this type of ES, and a class diagram is used to model facts by representing entities and their attributes. It is, therefore, establishing high-level specifications of business rules and processes. The generic PSM is developed based on pure production rules (FOPC), which makes it adaptable to different rule-based engines or Shells that implement PIM models of business aspects. Furthermore, we designed a UML profile diagram that extends the PIM sequence diagram, to support the lack of some features in the UML sequence model (OR and Not). Finally, in this tackle, we developed the necessary mapping rules (QVT) that act as a standard for the transformation of PIM sequence diagram metamodel into a rule-based PSM metamodel, generating the necessary rules and generating ES facts from UML class models as well as the developing round-trip mapping that supports the maintenance of ES.

To evaluate our proposed approach, which is design science research, a real case study of an academic advising system for low GPA students, was used for evaluation. QVT mapping rules that facilitate the transformation from the PIM to the PSM have been developed. In this process, we establish mapping rules that convert the PIM sequence diagram into a rule-based ES, generating the necessary ES rules. Additionally, we defined mapping rules that transform the PIM class diagram into a PSM rule-based ES, resulting in the creation of the required facts for ES. More importantly, utilizing the QVT Relational language that enables round-trip mapping thereby support potential changes (i.e. in rules, business requirement, platform) of PIM, PSM, mapping rules itself. A less costly maintenance therefore achieved because of the automation and the standardizing of round-trip mapping rules being developed. The results obtained from this case study provide practical evidence of the applicability and utility of our proposed approach in real-world scenarios.

Nevertheless, it is important to acknowledge the limitations of the current work. The connection between PSM and a platform is not tackled but since a generic PSM is developed the process is straightforward. Also, the consequences of the inclusion of OCL in UML models. In addition, although the

introduced model is adaptable to both process-oriented and data-oriented approaches, its primary focus lies in the process-oriented aspect. Also, the models lack the capability of using a relational or mathematical expression that can be needed in the PIM metamodel.

In future endeavors, our objective is to further advance the ES design approach by implementing and evaluating the proposed design on different domains of ES, ensuring its practical applicability and effectiveness, and supporting the lacking features in PIM. Also, incorporating the UML profile in the mapping process and resolving the limitation of tools (mapping engine) to recognize profiles.

REFERENCES

[1] H. Tan, "A brief history and technical review of the expert system research," IOP Conf. Ser. Mater. Sci. Eng., vol. 242, no. 1, 2017, doi: 10.1088/1757-899X/242/1/012111.

[2] R. Colomb, Deductive Databases and Their Applications. 1998.

[3] "The Oncology Expert Advisor," 2013. https://www.mdanderson.org/publications/annual-report/annual-report-2013/the-oncology-expert-advisor.html (accessed Dec. 08, 2023).

[4] B. T. Sayed, "Application of Expert Systems or Decision-Making Systems in the Field of Education," Inf. Technol. Ind., vol. 9, no. 1, pp. 1396–1405, 2021, doi: 10.17762/itii.v9i1.283.

[5] F. Lareyre, C. Adam, M. Carrier, and J. Raffort, "Automated segmentation of the human abdominal vascular system using a hybrid approach combining expert system and supervised deep learning," J. Clin. Med., vol. 10, no. 15, 2021, doi: 10.3390/jcm10153347.

[6] K. Fedra and L. Winkelbauer, "A hybrid expert system, GIS, and simulation modeling for environmental and technological risk management," Comput. Civ. Infrastruct. Eng., vol. 17, no. 2, pp. 131–146, 2002, doi: 10.1111/1467-8667.00261.

[7] A. M. Elsawi, S. Sahibuddin, and R. Ibrahim, "Model driven architecture a review of current literature," J. Theor. Appl. Inf. Technol., vol. 79, no. 1, pp. 122–127, 2015.

[8] S. Y. Choi and S. H. Kim, "Knowledge acquisition and representation for high-performance building design: A review for defining requirements for developing a design expert system," Sustain., vol. 13, no. 9, 2021, doi: 10.3390/su13094640.

[9] Y. Ran, X. Zhou, P. Lin, Y. Wen, and R. Deng, "A Survey of Predictive Maintenance: Systems, Purposes and Approaches," vol. XX, no. Xx, pp. 1–36, 2019, [Online]. Available: http://arxiv.org/abs/1912.07383.

[10] "OMG," 2014. https://www.omg.org/mda/ (accessed Feb. 02, 2022).

[11] "MDA," 2001. https://www.omg.org/mda/ (accessed Feb. 02, 2022).

[12] "ODM," Model Driven Eng. Ontol. Dev., no. September, pp. 215–233, 2009, doi: 10.1007/978-3-642-00282-3_8.

[13] "MOF," no. August, 2019, [Online]. Available: https://www.omg.org/spec/MOF/2.5.1/PDF.

[14] B. G. Buchanan and R. Q. Smith, "Fundamentals of expert system," Springer Ser. Mater. Sci., vol. 206, pp. 31–39, 1988, doi: 10.1007/978-3-662-44497-9_3.

[15] I. H. Sarker, A. I. Khan, Y. B. Abushark, and F. Alsolami, "Mobile expert system: Exploring context-aware machine learning rules for personalized decision-making in mobile applications," Symmetry (Basel)., vol. 13, no. 10, pp. 1–10, 2021, doi: 10.3390/sym13101975.

[16] N. Mayadevi, S. S. Vinodchandra, and S. Ushakumari, "A review on expert system applications in power plants," Int. J. Electr. Comput. Eng., vol. 4, no. 1, pp. 116–126, 2014, doi: 10.11591/ijece.v4i1.5025.

[17] S. S. A. Naser and M. H. Al-bayed, "Detecting Health Problems Related to Addiction of Video Game Playing Using an Expert System," J. Multidiscip. Res. Dev., vol. 2, no. 9, pp. 7–12, 2016.

[18] S. Khanna, A. Kaushik, and M. Barnela, "Expert Systems Advances in Education," Ncci, no. March, pp. 19–20, 2010, [Online]. Available: https://www.researchgate.net/profile/Akhil-Kaushik/publication/267862155.

[19] W. P. Wagner, "Trends in expert system development: A longitudinal content analysis of over thirty years of expert system case studies," Expert Syst. Appl., vol. 76, pp. 85–96, 2017, doi: 10.1016/j.eswa.2017.01.028.

[20] K. P. Tripathi, "A Review on Knowledge-based Expert System : Concept and Architecture," Artif. Intell. Tech. - Nov. Approaches Pract. Appl., vol. 4, no. 4, pp. 19–23, 2011.

[21] A. A. Mohammed, K. Ambak, A. M. Mosa, and D. Syamsunur, "Expert system in engineering transportation: A review," J. Eng. Sci. Technol., vol. 14, no. 1, pp. 229–252, 2019.

[22] C.-L. Chang and R. C. Tung lee, Symbolic Logic and Mechanical Theorem Proving. Academic Press, 1973.

[23] R. K. Lindsay, B. G. Buchanan, E. A. Feigenbaum, and J. Lederberg, "DENDRAL: A case study of the first expert system for scientific hypothesis formation," Artif. Intell., vol. 61, no. 2, pp. 209–261, 1993, doi: 10.1016/0004-3702(93)90068-M.

[24] "Drools." https://docs.drools.org/8.44.0.Final/drools-docs/drools/introduction/index.html (accessed Dec. 08, 2023).

[25] gensym, "G2." http://dev.gensym.com/platforms/g2-standard/# (accessed Dec. 10, 2023).

[26] F. Bellifemine, A. Poggi, and G. Rimassa, "JADE a FIPA2000 compliant agent development environment," Proc. Int. Conf. Auton. Agents, pp. 216–217, 2001.

[27] P. Charlton, R. Cattoni, A. Potrich, and E. Mamdani, "Evaluating the FIPA standards and their role in achieving cooperation in multi-agent systems," 2002, doi: 10.1109/HICSS.2000.926996.

[28] "OIDM," 2023. https://documentation.custhelp.com/euf/assets/devdocs/unversioned/Intelligent Advisor/en/Content/Guides/Overview/Overview.htm (accessed Dec. 12, 2023).

[29] J. Giarratano and G. Riley, Expert Systems: Principles and Programming, Fourth Edition. Course Technology, 2004.

[30] "QVT," Transformation, no. January, pp. 1–230, 2008, [Online]. Available: http://www.omg.org/spec/QVT/1.0/PDF/.

[31] V. Nikulsins, "Transformations of software process models to adopt model-driven architecture," Proc. 2nd Int. Work. Model. Archit. Model. Theory-Driven Dev. MDA MTDD 2010, Conjunction with ENASE 2010, pp. 70–79, 2010, doi: 10.5220/0003044500700079.

[32] David S. Frankel, Model Driven Architecture : Applying MDA to Enterprise Computing, vol. 308. 2003.

[33] R. S. Aguilar-Savén, "Business process modelling: Review and framework," Int. J. Prod. Econ., vol. 90, no. 2, pp. 129–149, 2004, doi: 10.1016/S0925-5273(03)00102-6.

[34] "UML," Proc. - 2005 IEEE Symp. Vis. Lang. Human-Centric Comput., vol. 2005, no. December, p. 9, 2005, doi: 10.1109/VLHCC.2005.65.

[35] "UML," pp. 443–506, 2017, [Online]. Available: https://www.omg.org/spec/UML/2.5.1/PDF.

[36] K. Jetlund, E. Onstein, and L. Huang, "Adapted rules for UML modelling of geospatial information for model-driven implementation as OWL ontologies," ISPRS Int. J. Geo-Information, vol. 8, no. 9, 2019, doi: 10.3390/ijgi8090365.

[37] "PRR," OMG Specif., vol. 1.0, no. December, p. 74, 2009, [Online]. Available: http://www.omg.org/spec/PRR/1.0/.

[38] I. Essebaa and S. Chantit, "Toward an automatic approach to get PIM level from CIM level using QVT rules," SITA 2016 - 11th Int. Conf. Intell. Syst. Theor. Appl., no. PMarch, 2016, doi: 10.1109/SITA.2016.7772271.

[39] R. S. Pressman and B. Maxim, Software Engineering: A Practitioner's Approach. McGraw Hill, 2010.

[40] "Jetpack Compose basics," 2023. https://developer.android.com/codelabs/jetpack-compose-basics#0 (accessed Aug. 05, 2023).

[41] "SwiftUI," 2023. https://developer.apple.com/documentation/swiftui/ (accessed Jul. 09, 2023).

[42] "WSDL," 2007. https://www.w3.org/TR/wsdl/ (accessed Jul. 05, 2023).

[43] "RESTful," 2018. https://wiki.onap.org/display/DW/RESTful+API+Design+Specification (accessed Dec. 12, 2023).

[44] "SOAP," 2007. https://www.w3.org/TR/2007/REC-soap12-part0-20070427/ (accessed Dec. 12, 2023).

[45] R. M. Tawafak, G. Alfarsi, A. Romli, J. Jabbar, S. I. Malik, and A. Alsideiri, "A Review Paper on Student-Graduate Advisory Expert system," 2020 Int. Conf. Comput. Inf. Technol. ICCIT 2020, pp. 187–191, 2020, doi: 10.1109/ICCIT-144147971.2020.9213794.

[46] "EMF Tutorial - EclipseSource." https://eclipsesource.com/blogs/tutorials/emf-tutorial/ (accessed Oct. 18, 2022).

# Comparison of SVM kernels in Credit Card Fraud Detection using GANs

Bandar Alshawi

Department of Computer and Network Engineering, College of Computing, Umm Al-Qura University, Makkah, Saudi Arabia

*Abstract*—**The technological evolution in smartphones and telecommunication systems have led people to be more dependent on online shopping and electronic payments, which created burdensome task of transaction validation for many financial institutions. This paper examined and evaluated the efficacy of Support vector machine (SVM) kernels on Generative Adversarial Network (GAN)-generated synthetic data to detect credit card fraud transactions. Four SVM kernels have been investigated and compared; linear, polynomial, sigmoid, and redial basis function. The accuracy results indicated that linear and polynomial kernels reached over 91%, while sigmoid and redial basis function reached 79% and 83% respectively. Linear and polynomial models received over 90% ROC and F1 score, in contrast the ROC scores were lower for sigmoid (81%) and redial basis function (83%). Both sigmoid and redial basis function achieved over 80% in terms of F1 score. The precision score demonstrated a high score for both linear and polynomial kernel reaching 99%. Additionally, sigmoid and redial basis function achieved over 80%. These results overcame the imbalance dataset issue through the generation of synthetic data by applying the SVM kernels using GANs algorithm.**

*Keywords*—*Fraud transactions; credit card; Generative Adversarial Network; Support Vector Machine kernels; imbalance dataset*

## I. INTRODUCTION

The evolution of telecommunications technologies and the adoption of electronic payments from vast financial institutions led to unanticipated spike in fraud transactions. Personal and organizational assets nowadays are vulnerable due to cybersecurity breaches [1]. In 2020 alone, banks have suffered over $28 billion in credit card losses globally. The numbers are predicted to surpass $49 billion by 2030 [2]. Engaging artificial intelligence in banking system will enhance fraud detection, thus protecting assets and reinforce customer fidelity [3]. The fraud and control report in [4] sheds light to almost 26% of electronic transactions were categorized as fraud or attempted fraud. Detecting electronic fraud transactions using Machine Learning (ML) can be cumbersome according to the research presented in [5]. Diverse ML credit card fraud detection system has been previously reviewed [6, 7]. The complexity of imbalance dataset exists in different real-world ML scenarios. In the credit card dataset, the irregular distribution of one class was evident due to the fact that valid transaction exceeds fraudulent transactions [8]. Numerous credit card fraud detection methods capable of avoiding fraudulent transactions in the banking sectors include data mining, modeling algorithms, which comprise of clustering methods and fraud detection [9].

This research investigates an important issue in credit card fraud detection using ML techniques, which raises the following questions:

- Has any of the previous research examined different SVM kernels to detect credit card fraud transactions on imbalanced dataset?

- How does the Generative Adversarial Network (GAN) perform on generating tabular data?

- How the four SVM kernels perform against each other.

To answer the preceding questions of this research, numerous objectives required to be met, including:

- Reparation of imbalance dataset in tabular data.

- Using specific GAN to generate synthetic tabular data.

- Detecting credit card fraud using SVM kernels and evaluate the performance of each kernel among other.

Although GANs are mainly used to synthesize visual data, several research have successfully managed to use them to generate tabular. The significance of this study is overcoming the issue of imbalanced dataset while investigating the performance of different SVM kernels.

This paper is categorized as follows: Section II presents related efforts on several ML fraud detection study; Section III discusses the methods used to predict the results. In Section IV an extensive review of the results and analysis is detailed. Section V is reserved for discussion and comparisons. Section VI presents the conclusion of the research.

## II. RELATED WORK

The research in study [10] presented an approach to observe credit card fraud. The author focused on reaching unbiased and consistent techniques to automate fraud risk evaluation. The approach proposed an algorithm that calculated variables' relationships and related information. The solution successfully improved accuracy and diminished dimensionality. The study in [11] illustrated a comparison of various credit card fraud detection methods using supervised and unsupervised learning. The results show a prime for unsupervised learning, while emphasizing the effects on performance when using supervised learning methods. A fraud financial detection method was presented in [1] named Intimation Rule Based (IRB) alert generation algorithm using ontology-based system which benefited from ontology alert. The author constructed their method by including forty categories and sub-categories which effectively can capture

fraud by sending different notifications according to their extremity. In study [12] the author examined the utilization of supervised and unsupervised methods to identify inconsistencies in financial transaction records. The research in [13] proposed a hybrid ensemble model to detect anomalies in credit card transactions. The research used adaboost, random forest, and logistic regression as classifiers, imbalanced dataset was addressed by oversampling method and removal of outliers. The study examined SVM along with different ML methods including an adaptive boosting (AdaBoost) and decision tree on real world dataset. The experiment involved the use of real-world dataset and applying vectorization on the sub-leader account size to tangle irregularity. Most classifiers are incapable of procuring acceptable outcome during imbalance data classification, the author in [14] proposed an optimized SVM by Genetic Algorithm, dataset balancing is done through cluster centroids sampling.

The study in [15] applied SVM along with decision trees on an extremely imbalanced real-world dataset. A handful of numbers of machine learning techniques were examined that include outlier detection and ensemble algorithms. The author employed feature engineering to calculate the effect of feature-selection on performance. The research in [16] reviewed the latest progress in detecting fraud transaction using Deep Reinforcement Learning (DRL) and ML. The research carried out an experiment on an exceedingly imbalanced dataset using resampling technique to deal with complications and implementing several ML and DRL methods. An extensive analysis was carried out on non-linear models in [17]. The study proposed binary types of fraud detector models, one that can be interpreted and the other cannot be bound to a specific way. The models are utilized concurrently with ML methods. Furthermore, Black Box model is avoided in the study by supply tracing information that associates inputs and outputs. Credit card fraud detection methods using several neural networks concurrently with resembling methods were demonstrated in [18]. A combination of Harris Hawk Optimization (HHO) and SMOTE was introduced in [19]. The study tried to identify the appropriate sampling pace for the HHO and combines it with the SMOTE algorithm. The main aim of the study is to maximize classifier accuracy in imbalanced datasets. Different ML techniques were discussed such as: recurrent neural network, convolutional neural network, and ensemble methods. The research attempted to investigate obstacles and limitations related to IoT anomaly detectors.

The research in study [20] presented a system that integrates Deep Neural Network (DNN) and Catboost, to test any overlap in classification rate improvement. The experiment was carried out on IEEE-CIS dataset composed of 590,540 instances. Miscellaneous classifiers have been tested on highly imbalanced datasets in [21]. The author applied random over sampling (RO), which replicate instances from the minority category followed by applying SVM, NÏVE BAYS (NB), Artificial neural network (ANN), and C5.0. The review in [22] discussed oversampling and undersampling to handle imbalance dataset and comparing convolutional to an ensemble algorithm during credit card fraud detection, concluding that ensemble was more effective. An ensemble methodology was used by applying decision tree, logistic regression, and NB side by side in [23] highest output is picked by hard voting. In study [24] rough set theory was used for initial data refinement consisting of attribute estimation and reduction, lease square support vector later applied to classify and predict credit card churn behavior. Hierarchical temporal memory, based on cortical learning HTM-CLA algorithm, was presented in [25] to recognize fraudulent transaction. The authors also measured the difference between the HTM-CLA outcomes of using traditional Artificial Neural Network tree (ANN) in contrast to simulated annealing ANN. The research in study [26] used GANs along with logistic regression, decision tree, naïve bay, random forest, extreme gradient boosting, and adaptive boosting algorithms to detect credit card fraud transactions.

## III. METHODOLOGIES

### A. Implementation

The used tools through this research include intel i9-9900K 3.60GHz, 64GB RAM, Nvidia 2080TI was utilized for GAN synthetic data generation and SVM kernels training and testing.

### B. Original Data

The dataset contains credit card transactions by European cardholders in October 2013. The dataset consisted of transactions that occurred in two consecutive days. The dataset is imbalanced since it had 492 flagged as fraudulent out of the 284,807 transactions. With accordance to client's confidentiality and privacy, the dataset underwent the Principal Component Analysis (PCA), resulting in numerical variables. The dataset consisted of 31 features that are Class, Time, V1, and V28.

### C. Synthetic Data

Different researchers have tried to cope with the complexity of imbalanced dataset in the existing area using Synthetic Minority Oversampling Technique (MOTE) found in [27], [12], and [28]. SMOTE is an effective oversampling technique used to generate synthetic data from minority class [29]. Synthetic data was adopted in [30] using Monte Carlo simulations, a whole dataset was assembled including a number of features. GANs have been adopted in numerous domains recently to refine synthetic images producing a realistic representation. Other example of GANs adoption was done by Alonso et al. in [31] their model generates handwritten text. Their generator is conditioned on a sequence of characters, subsequently the generator starts producing synthetic data in the form of handwritten instances for different words. Creating synthetic data out of a random noise is not GANs prominent purpose, its capability lies in estimating the uneven class and generating data from a small set of samples [32]. The proposed approach uses synthetic data, which can be generated using GANs method discussed and explained thoroughly in Section III.

### D. Generative Adversarial Network

GAN was created by Ian Goodfellow in 2014, and it is classified under unsupervised learning [33]. Generative models are capable of learning and imitating any distribution of data. GANs consist of two neural networks trained competitively thus; they are referred to as adversarial. GANs utilize the deep

neural network as a training algorithm. Imbalanced dataset is a frequent matter during modeling and may result in a weak model therefore; GANs can be employed to generate synthetic data which could solve some of the complexity [34]. GANs consist of binary neural networks operating in a contrary mode, the former is identified as the generator, and the latter refers to as the discriminator. Collection and generation of samples are the purpose of Generator Network presence. The probability of discriminator network to mis-classify would grow proportionally during the training of the generator network. GANs equation is found in Eq. (1) where G is the generative

model learning from the training data x, D is the discriminator, which separate among various classes of data. The discriminator identifies whether the received data were generated from a real sample using a binary for output ranging from 0 to 1. In Eq. (1) the generator receives a slight noise sample from z. ~µ_z refers to the generator distribution and ~µ_ref refers to the real data distribution. GANs architecture is presented in Fig. 1.

$$L(G, D) = \mathbb{E}_{x \sim \mu_{ref}}\left[\ln D(x)\right] + \mathbb{E}_{z \sim \mu_z}\left[ln\left(1 - D\big(G(Z)\big)\right)\right]$$
(1)



Fig. 1.    Generative adversarial network.

### E. GANs for Tabular Data

Applying GANs over tabular data, can rise numerous of challenges such of which are indicated below:

- Tabular data can be of a mixed type.

- GANs are effective in image data, and they distribute them over space. On the other hand, tabular data are non-Gaussian that could affect the network not being able to propagate gradient details.

- The generator is not capable of recognizing imbalanced categorical columns when using generated samples from standard multivariate distribution.

CTGAN is implemented in this research which is a GANs based method capable of solving non-Gaussian obstacle by applying mode-specific normalization. Moreover, it uses all existing features of the dataset [35]. In Fig. 2 an illustration of the proposed framework of credit card detection approach is demonstrated.

### F. Machine Learning

ML was defined since 1959 by the AI pioneer Arthur Samuel who indicated that computers will be capable of learning from experiences rather than being programmed. ML is classified into three categories: supervised learning, unsupervised learning, and reinforcement machine learning [36]. Fig. 3 illustrates SVM kernel classification utilized in this research.



Fig. 2.    Framework of the proposed credit card detection approach.

Fig. 3.    SVM kernel classification.

### G. Support Vector Machine

The SVM algorithm is a supervised learning algorithm developed in Bell Lab by Vladimir Vapnik. SVM has the capability of solving regression and classification problems. SVM can separate two categories by drawing a hyperplane. The performance of SVM is thoroughly impacted by the selected kernel and the parameters as default or set values. The use of kernel aims to assemble a nonlinear hyperplane including all the set of input values to execute the classification [37].

### H. Kernel Functions

Kernel functions come to place in situations where samples are linearly non-sparable. SVM kernel includes decision functions to non-linear class by mapping input sample and projecting them into a higher dimensional space, not requiring calculating the mapping explicitly. Optimistically, the samples will achieve significant linear structure. Moreover, the kernel function can be thought of as a measure of similarity between samples [38], which grants SVM to carry out separations regardless of very complex boundaries. Different kernel settings will be discussed the following section.

### I. Radial Basis Function RBF Kernel

Radial Basis Function falls under the neural network types. It is used to solve diverse problems such as: classification, prediction, and regression. The approach that RBF uses to process classification problems differ compared to ordinary neural network. Ordinary neural network performs data separation using linear manipulations of activation function. On the other hand RBF organize data by density-based transformation. RBF equation is stated below where X  is the input,  C is the mean center between lowering the training error and surging margin[15], and σ is the spread.

$$h(X) = exp\left(-\frac{\|X - C\|^2}{2\sigma^2}\right) \tag{2}$$

### J. Polynomial Kernel

Polynomial kernel is another type of SVM kernel, while it benefits from the polynomial function. The polynomial kernel is used to resale data into greater space. This operation is done by taking the scalar product of data points, the existing space with the polynomial in the newer space. Using polynomial functions allows for greater dimensional mapping for data. The equation of polynomial kernel is shown below where $x$ and $y$ are vectors and $c$ is a constant in the existing space. $d$ is the degree of the polynomial function.

$$K(x,y) = (x^T y + c)^d \tag{3}$$

### K. Sigmoid Kernel

The idea of sigmoid kernel evolved from neural networks. Sigmoid kernel usage can be problematic due parameters adjustments [39]. The equation of Sigmoid kernel is shown below where $X_i$ and $X_j$ are vectors, $\beta_0$ is the slope, and $\beta_1$ is the intercept.

$$K(X_i, X_j) = tanh(\beta_0 X_i^T X_j + \beta_1) \tag{4}$$

### L. Linear Kernel

Linear kernel is the simplest kernel, unlike the polynomial or the logistic regression where data are projected to the upper space. In linear kernel, it obtains a single dimensional nature. In other words, linear kernel is capable of separating classes using the hyperplane with linear boundaries [40]. The equation of linear kernel is listed at Eq. (5) where $X, X_j$ represents the data to be classified.

$$K(X, X_j) = sum(X \cdot X_j) \tag{5}$$

### M. Evaluation Metrics of SVM Kernels

SVM kernels performance were evaluated and tested using confusion matrix. Which contain True Positive (TP) to correspond to valid transactions that were predicted correctly, false positive (FP) refers to fraud transactions that were not captured, true negative (TN) indicate fraud transactions that were predicted accurately, and false negative (FN) illustrating fraud transactions that were not identified as fraud by the model. The equation presented in Eq. (6) was used to measure accuracy of correctly predicted transactions over the total number of transactions, sometimes referred to as error rate.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \tag{6}$$

Precision metric computes positive instances that are accurately predicted from the total positive predictions. The equation of precision is presented in Eq. (7). One thing to note is that precision and recall do operate in contrast usually one is higher than the other.

$$Precision = \frac{TP}{(TP+FP)} \tag{7}$$

Recall metric displays a calculation of a portion of positive instances that are accurately classified. Their importance stems from their capability of capturing positive cases and that higher recall value prevents missing fraudulent transaction [3]. The equation of recall is found in Eq. 8.

$$Recall = \frac{TP}{(TP+FN)} \tag{8}$$

F1 metric is used to obtain the harmonic mean between precision and recall. F1 has a score between 0 and 1, higher values reflect high model performance. The equation of F1 appears at Eq. (9).

$$F - measure = \frac{2*Precision*Recall}{Precision+Recall} \tag{9}$$

Matthews Correlation Coefficient (MCC) was invented by Brian Matthews in 1975. MCC measures the quality of the classifiers between observation and prediction. It can be described as a confusion matrix method of calculating the Pearson

product-moment correlation coefficient between predicted and actual value [16].

$$MCC = \frac{TP * TN - FP * FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (10)$$

## IV. RESULTS

Metric evaluation is summarized in Table I and Table II. An additional method to evaluate the results is through the confusion matrix, which can be used for both binary classification as well as multiclass classification. In the current state, which is binary classification, the confusion matrix generates table of 2*2. The output consists of TP, FP, TN, and FN which were discussed earlier in section III. This section illustrates the outcome of the four SVM kernels discussed in Section III. The accuracy results in Table I showed that LN kernel score was the highest at 95% followed by PL at 91%. RBF scored 83% and SG scored 79%. In precision both LN and PL scored 99%, followed by RBF at 84% and lastly SG which was at 81%. The recall results demonstrated PL with 85% success rate by RBF 85%, led by LN at 82% and SG at 81%. In Table II, the Receiver Operating Characteristic Curve (ROC) is utilized to calculate classifier performance at different thresholds. ROC score in Table II display PL and LN achieved 91%, while RBF attained 83% and SG reached 79%. Based on the MCC score, LN achieved 84%, followed by PL at 83% succeeded by RBF at 66%, and lastly SG at 58%. F1 score in Table II indicated the harmonic mean between precision and recall. LN has achieved the highest F1 score at 92% followed by PL at 91%, concluded with RBF at 85% and SG at 81%. Confusion matrix is used in this section to ease the understanding of outcome, which generates 2*2 table with binary values representing multiclass classification; TP, TN, FP, and FN discussed in Section III. Fig. 4 presents the confusion matrix for LN kernel indicating that the classifier accurately predicted 4285 legitimate and 4145 fraud transactions. LN kernel inaccurately classified 74 fraud transactions as legitimate and 715 legitimate transactions as fraud. Fig. 5 displays the confusion matrix for PL kernel point that the classifier accurately predicted 4253 valid and 4160 invalid transactions. PL classifier was not able to capture 62 fraud transactions furthermore classified 747 real transaction as fraud. Fig. 6 shows the confusion matrix for RBF kernel, which predicted 4231 legitimate and 3468 fraud transactions accurately. Fig. 6 also exposes that 781 fraudulent transactions were classified as valid, and 769 valid transactions were categorized as fraud. SG confusion matrix appears at Fig. 7, representing that SG kernel correctly classified 4027 legitimate and 3305 fraud transactions. SG kernel was not able to capture 938 fraud transactions and categorized them incorrectly; it also categorized 973 valid transactions as fraud. Fig. 8 illustrates the AUC – ROC curve, which explains the classification performance at different thresholds. From the figure it is evident that PL and LN classifiers have a better measure of separability than SG and RBF classifiers. The precision-recall curve is presented in Fig. 9, which can be employed especially when imbalance dataset is in existence. As the figure illustrates PL and LN did outperform SG and RBF.

TABLE I. PERFORMANCE EVALUATION

| Algorithm | Accuracy | Precision | Recall |
|---|---|---|---|
| Polynomial | 91 | 99 | 85 |
| Sigmoid | 79 | 81 | 81 |
| Linear | 95 | 99 | 82 |
| RBF | 83 | 84 | 85 |

TABLE II. PERFORMANCE EVALUATION: ROC, MCC, AND F1-SCORE

| Algorithm | ROC | MCC | F1 score |
|---|---|---|---|
| Polynomial | 91 | 83 | 91 |
| Sigmoid | 79 | 58 | 81 |
| Linear | 91 | 84 | 92 |
| RBF | 83 | 66 | 85 |



Fig. 4. Linear kernel confusion matrix.



Fig. 5. Polynomial kernel confusion matrix.

Fig. 6.    RBF kernel confusion matrix.



Fig. 7.    Sigmoid kernel confusion matrix.



Fig. 8.    ROC.



Fig. 9.    PRC.

## V.    DISCUSSION

In Table III a summary of relevant studies is presented, starting with the research in [41], which examined two kernels LN and RBF. Based on the results from Table III it is evident that there is a huge variation between accuracy and recall. The research in [21, 42, 43, 44] does not investigate SVM kernels. In [42, 44] both studies evaluated their models with accuracy which is not always an accurate metric indicator. A comparative result is demonstrated in Fig. 10 which presents consistent overall results for the proposed solution compared to recent existing studies. Few research examined the detection of credit card fraud detection using different SVM kernels, none of them produced significant outcome as noted in Table III and Fig. 10.

TABLE III.    COMPARISON OF RELEVANT PAST STUDIES

| | Classifier | Kernel | Accuracy (%) | Recall (%) | Year of publication | Reference |
|---|---|---|---|---|---|---|
| 1 | SVM | - | 96.34 | - | 2022 | [44] |
| 2 | SVM | - | 99 | - | 2022 | [43] |
| 3 | SVM | - | 96 | 39 | 2019 | [21] |
| 4 | SVM | - | 99 | - | 2018 | [42] |
| 5 | SVM | LN | 97 | 1 | 2020 | [41] |
| | | | 97 | 10 | | |
| | | RBF | 97 | 86 | | |
| | | | 97 | 82 | | |

Fig. 10. Comparative results with relevant studies.

## VI. CONCLUSION

The recent advancements and improvements in technological and telecommunication industry lead various sectors to integrate this technology into their system. In addition, due to the tremendous increase of electronic transactions, financial institutions are affected by fraudulent transactions, which meant that certain procedures must take place, including the adoption of ML fraud prevention techniques. In this paper, GAN was used to generate synthetic data to overcome uneven class distribution of credit card dataset. Four SVM kernels were used to predict fraudulent transactions and compared with each other and with relevant recent research. The findings illustrated that two SVM kernels LN and PL scored over 91% in accuracy however, RBF achieved 83% while SG reached 79%. LN and PL have received an over 91% ROC and F1 scores, yet SG reached 79% and RBF scored 83% in ROC. The F1 score for SG and RBF demonstrate that both kernels received over 81%. The future work should focus on investigating the use of different GAN variants with SVM and different classifiers.

## REFERENCES

[1] M. Ahmed, K. Ansar, C. B. Muckley, A. Khan, A. Anjum, and M. Talha, "A semantic rule based digital fraud detection," *PeerJ Computer Science,* vol. 7, p. e649, 2021.

[2] D. Robertson, "Card Fraud Worldwide," Nelson report, 2021.

[3] N. S. Alfaiz and S. M. Fati, "Enhanced credit card fraud detection model using machine learning," *Electronics,* vol. 11, no. 4, p. 662, 2022.

[4] J. P. Morgan, "Payments Fraud and Control Report," 2022.

[5] A. Dal Pozzolo, G. Boracchi, O. Caelen, C. Alippi, and G. Bontempi, "Credit card fraud detection: a realistic modeling and a novel learning strategy," *IEEE transactions on neural networks and learning systems,* vol. 29, no. 8, pp. 3784-3797, 2017.

[6] S. N. Kalid, K.-H. Ng, G.-K. Tong, and K.-C. Khor, "A multiple classifiers system for anomaly detection in credit card data with unbalanced and overlapped classes," *IEEE access,* vol. 8, pp. 28210-28221, 2020.

[7] M. C. M. Oo and T. Thein, "An efficient predictive analytics system for high dimensional big data," *Journal of King Saud University-Computer and Information Sciences,* vol. 34, no. 1, pp. 1521-1532, 2022.

[8] A. Dal Pozzolo, O. Caelen, R. A. Johnson, and G. Bontempi, "Calibrating probability with undersampling for unbalanced classification," in *2015 IEEE symposium series on computational intelligence*, 2015: IEEE, pp. 159-166.

[9] D. Dighe, S. Patil, and S. Kokate, "Detection of credit card fraud transactions using machine learning algorithms and neural networks: A comparative study," in *2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA)*, 2018: IEEE, pp. 1-6.

[10] J. Chaquet-Ulldemolins, F.-J. Gimeno-Blanes, S. Moral-Rubio, S. Muñoz-Romero, and J.-L. Rojo-Álvarez, "On the black-box challenge for fraud detection using machine learning (I): Linear models and informative feature selection," *Applied Sciences,* vol. 12, no. 7, p. 3328, 2022.

[11] X. Niu, L. Wang, and X. Yang, "A comparison study of credit card fraud detection: Supervised versus unsupervised," *arXiv preprint arXiv:1904.10604,* 2019.

[12] A. Bakumenko and A. Elragal, "Detecting anomalies in financial data using machine learning algorithms," *Systems,* vol. 10, no. 5, p. 130, 2022.

[13] S. Saraf and A. Phakatkar, "Detection of Credit Card Fraud using a Hybrid Ensemble Model," *International Journal of Advanced Computer Science and Applications,* vol. 13, no. 9, 2022.

[14] Y. Cui, Z. Song, and J. Hu, "esearch on Credit Card Fraud Classification Based on GA-SVM," in *2021 4th International Conference on Advanced Electronic Materials, Computers and Software Engineering (AEMCSE)*, 2021: IEEE, pp. 1076-1080.

[15] Y. G. Şahin and E. Duman, "Detecting credit card fraud by decision trees and support vector machines," 2011.

[16] T. K. Dang, T. C. Tran, L. M. Tuan, and M. V. Tiep, "Machine learning based on resampling approaches and deep reinforcement learning for credit card fraud detection systems," *Applied Sciences,* vol. 11, no. 21, p. 10004, 2021.

[17] J. Chaquet-Ulldemolins, F.-J. Gimeno-Blanes, S. Moral-Rubio, S. Muñoz-Romero, and J.-L. Rojo-Álvarez, "On the black-box challenge for fraud detection using machine learning (ii): nonlinear analysis through interpretable autoencoders," *Applied Sciences,* vol. 12, no. 8, p. 3856, 2022.

[18] E. Esenogho, I. D. Mienye, T. G. Swart, K. Aruleba, and G. Obaido, "A neural network ensemble with feature engineering for improved credit card fraud detection," *IEEE Access,* vol. 10, pp. 16400-16407, 2022.

[19] K. S. Raslan, A. S. Alsharkawy, and K. R. Raslan, "HHO-SMOTe: Efficient Sampling Rate for Synthetic Minority Oversampling Technique Based on Harris Hawk Optimization," *International Journal of Advanced Computer Science and Applications,* 2023.

[20] N. Nguyen *et al.*, "A proposed model for card fraud detection based on Catboost and deep neural network," *IEEE Access,* vol. 10, pp. 96852-96861, 2022.

[21] S. Makki, Z. Assaghir, Y. Taher, R. Haque, M.-S. Hacid, and H. Zeineddine, "An experimental study with imbalanced classification approaches for credit card fraud detection," *IEEE Access,* vol. 7, pp. 93010-93022, 2019.

[22] A. N. Ahmed and R. Saini, "A Survey on Detection of Fraudulent Credit Card Transactions Using Machine Learning Algorithms," in *2023 3rd International Conference on Intelligent Communication and Computational Techniques (ICCT)*, 2023: IEEE, pp. 1-5.

[23] P. Tomar, S. Shrivastava, and U. Thakar, "Ensemble Learning based Credit Card Fraud Detection System," in *2021 5th Conference on Information and Communication Technology (CICT)*, 2021: IEEE, pp. 1-5.

[24] N. Wang and D.-x. Niu, "Credit card customer churn prediction based on the RST and LS-SVM," in *2009 6th international conference on service systems and service management*, 2009: IEEE, pp. 275-279.

[25] E. Osegi and E. Jumbo, "Comparative analysis of credit card fraud detection in Simulated Annealing trained Artificial Neural Network and Hierarchical Temporal Memory," *Machine Learning with Applications,* vol. 6, p. 100080, 2021.

[26] B. Alshawi, "Utilizing GANs for Credit Card Fraud Detection: A Comparison of Supervised Learning Algorithms," *Engineering, Technology & Applied Science Research,* vol. 13, no. 6, pp. 12264-12270, 2023.

[27] E. Ileberi, Y. Sun, and Z. Wang, "Performance evaluation of machine learning methods for credit card fraud detection using SMOTE and AdaBoost," *IEEE Access,* vol. 9, pp. 165286-165294, 2021.

[28] Z. Li, G. Liu, and C. Jiang, "Deep representation learning with full center loss for credit card fraud detection," *IEEE Transactions on Computational Social Systems,* vol. 7, no. 2, pp. 569-579, 2020.

[29] T. C. Tran and T. K. Dang, "Machine learning for prediction of imbalanced data: Credit fraud detection," in *2021 15th International Conference on Ubiquitous Information Management and Communication (IMCOM)*, 2021: IEEE, pp. 1-7.

[30] A. Singh, J. Amutha, J. Nagar, S. Sharma, and C.-C. Lee, "AutoML-ID: Automated machine learning model for intrusion detection using wireless sensor network," *Scientific Reports,* vol. 12, no. 1, p. 9074, 2022.

[31] E. Alonso, B. Moysset, and R. Messina, "Adversarial generation of handwritten text images conditioned on sequences," in *2019 international conference on document analysis and recognition (ICDAR)*, 2019: IEEE, pp. 481-486.

[32] S. I. Nikolenko, "Synthetic data for deep learning," *arXiv preprint arXiv:1909.11512,* 2019.

[33] I. Goodfellow *et al.*, "Generative adversarial nets," *Advances in neural information processing systems,* vol. 27, 2014.

[34] I. Ashrapov, "Tabular GANs for uneven distribution," *arXiv preprint arXiv:2010.00638,* 2020.

[35] L. Xu, M. Skoularidou, A. Cuesta-Infante, and K. Veeramachaneni, "Modeling tabular data using conditional gan," *Advances in neural information processing systems,* vol. 32, 2019.

[36] S. B. Kotsiantis, I. D. Zaharakis, and P. E. Pintelas, "Machine learning: a review of classification and combining techniques," *Artificial Intelligence Review,* vol. 26, pp. 159-190, 2006.

[37] T. M. T. Ab Hamid, R. Sallehuddin, Z. M. Yunos, and A. Ali, "Ensemble based filter feature selection with harmonize particle swarm optimization and support vector machine for optimal cancer classification," *Machine Learning with Applications,* vol. 5, p. 100054, 2021.

[38] R. Amami, D. B. Ayed, and N. Ellouze, "Practical selection of SVM supervised parameters with different feature representations for vowel recognition," *arXiv preprint arXiv:1507.06020,* 2015.

[39] H.-T. Lin and C.-J. Lin, "A study on sigmoid kernels for SVM and the training of non-PSD kernels by SMO-type methods," *Neural Comput,* vol. 3, no. 1-32, p. 16, 2003.

[40] C. Savas and F. Dovis, "Comparative performance study of linear and gaussian kernel SVM implementations for phase scintillation detection," in *2019 International Conference on Localization and GNSS (ICL-GNSS)*, 2019: IEEE, pp. 1-6.

[41] A. A. Taha and S. J. Malebary, "An intelligent approach to credit card fraud detection using an optimized light gradient boosting machine," *IEEE Access,* vol. 8, pp. 25579-25587, 2020.

[42] K. Randhawa, C. K. Loo, M. Seera, C. P. Lim, and A. K. Nandi, "Credit card fraud detection using AdaBoost and majority voting," *IEEE access,* vol. 6, pp. 14277-14284, 2018.

[43] S. Khan, A. Alourani, B. Mishra, A. Ali, and M. Kamal, "Developing a Credit Card Fraud Detection Model using Machine Learning Approaches," *International Journal of Advanced Computer Science and Applications,* vol. 13, no. 3, 2022.

[44] F. K. Alarfaj, I. Malik, H. U. Khan, N. Almusallam, M. Ramzan, and M. Ahmed, "Credit card fraud detection using state-of-the-art machine learning and deep learning algorithms," *IEEE Access,* vol. 10, pp. 39700-39715, 2022.

# A Cost-Efficient Approach for Creating Virtual Fitting Room using Generative Adversarial Networks (GANs)

Kirolos Attallah[1], Girgis Zaky[2], Nourhan Abdelrhim[3], Kyrillos Botros[4], Amjad Dife[5], Nermin Negied[6]

Faculty of Electrical and Computer Engineering, University of Ottawa, Ottawa, Canada[1, 2, 3, 4, 5]

School of Information Technology and Computer Science, Nile University, Giza, Egypt[6]

*Abstract*—**Customers all over the world want to see how the clothes fit them or not before purchasing. Therefore, customers by nature prefer brick-and-mortar clothes shopping so they can try on products before purchasing them. But after the Pandemic of COVID19 many sellers either shifted to online shopping or closed their fitting rooms which made the shopping process hesitant and doubtful. The fact that the clothes may not be suitable for their buyers after purchase led us to think about using new AI technologies to create an online platform or a virtual fitting room (VFR) in the form of a mobile application and a deployed model using a webpage that can be embedded later to any online store where they can try on any number of cloth items without physically trying them. Besides, it will save much searching time for their needs. Furthermore, it will reduce the crowding and headache in the physical shops by applying the same technology using a special type of mirror that will enable customers to try on faster. On the other hand, from business owners' perspective, this project will highly increase their online sales, besides, it will save the quality of the products by avoiding physical trials issues. The main approach used in this work is applying Generative Adversarial Networks (GANs) combined with image processing techniques to generate one output image from two input images which are the person image and the cloth image. This work achieved results that outperformed the state-of-the-art approaches found in literature.**

*Keywords—Generative Adversarial Networks (GANs), virtual reality; human body segmentation; image generator; conditional generator; background removal*

## I. INTRODUCTION

Virtual fitting rooms (VFR) bring great opportunities to the fashion industry by enabling consumers to virtually try on products. However, while VFRs have technically been available for a while, they are less utilized because of many reasons, amongst them the consumers' potential concerns of accuracy of the simulation, the cost of the technologies used in building VFRs like Kinect cameras and depth sensors, and the difficulty of using them by the customer because of the special settings they require like special types of LCDs and mirrors. Research has proven that online clothes purchases had increased, and the return requests had decreased after the intervention of virtual fitting rooms (VFRs) [1]. e-Commerce development, AI development, and pandemics shared in both the research and industry demand of VFRs. In this paper two novel outputs are obtained and explained, the first one is a new dataset which is large and various (i.e., containing images for males and females in different poses and from different fashion houses, without excluding any challenges from the collected images) for the purpose of building VFRs, and the second one is a real time, cost efficient, portable, easy to use, and accurate VFR. The rest of this paper is organized as follows: Section II reviews the work done in literature to address this problem. Section III explains the methodologies and techniques used in this work to build the VFR. Section IV demonstrates the experiments conducted to evaluate the work. Section V discusses the results. Finally, the paper is concluded in Section VI.

## II. LITERATURE REVIEW

Recently, the idea of virtual fitting rooms has attracted researchers because of the emergence and development of virtual and augmented reality. The digital revolution and affordable technologies and devices also made virtual fitting rooms an area of interest. According to literature there are eight different types of virtual fitting rooms (VFR), which are: Body scanning VFRs, 3D avatar VFRs, 3D customer's modelling VFR, 3D mannequin VFR, Augmented reality VFRs, Robotic mannequins, Dress-up mannequins for mix-and-match, and the real fashion model VFR [1 & 2]. From all technologies and types of body scanning machines proved to be the most accurate, but the most expensive at the same time [3-6]. Although the large variety of technologies that could share in maximizing the customer experience while trying the cloth item online, the cost plays a major role of the types existing in the market. Some features and complementary methods could share in maximizing the customer experience besides the VFRs such as: fit guides, size charts, comparison avatars, virtual cat walks, brands comparisons, etc. [7 - 8].

Researchers used many emergent technologies to simulate fitting rooms, like augmented reality (AR), virtual reality (VR), depth information using Kinect cameras, deep learning approaches (DL), 3D reconstruction, and hybrid approaches. Pereira et al. in 2010 [9], used AR and Kinect camera to establish virtual fitting room using depth information, the authors tested their approach using Open CV and Open GL techniques with different six degrees of human head detection, but they mentioned nothing factual about their results.

Kostas [10] used the same technologies (Kinect cameras and AR) to implement virtual fitting room for trying different cloth items, but the higher accuracy they reached was 67%. Dias et al. [11] also used the same approaches for the same purpose, but they added different features, and they developed a better user interface, but they mentioned nothing about the VFR results. Below is a sample of their results. Hashmi et al. [12] used the AR combined with Haar-cascades classifier to implement VFR. The authors tested their classifier using 50

subjects with 10 different dresses and they confirmed that their approach had outperformed other approaches in literature. Mehta et al. [13] had used AR, VR, and Mixed reality (MR) combined with Head Mounted Displays (HMD). The authors claimed that their VFR can give customers better online shopping experience, but again they mentioned nothing factual about the results.

Boonbrahm et al. [14] in 2015 used VR to build VFR. The authors considered the differences between materials of clothes, and they confirmed that this matters in the final appearance of the dressed person. The materials they considered were jean, satin, silk, and cotton. França & Soares [15] in 2018 have discussed the idea of a complete simulation of Virtual Fitting Room (VFR) using virtual reality (VR) in which the customer feels he/she is existing in a real fitting room. Alfredo et al. [16] used Kinect camera along with gesture recognition and cloth transfer algorithm to create a virtual try in application. Despite the authors concluded that their approach combined with the depth information obtained using the Kinect camera made the customer experience more enjoyable, they mentioned nothing factual about the results.

Sapio et al. [17] in 2018 integrated different body scanning technologies to create a VFR. The authors used Kinect combined with other more expensive body scanners, but they found out that their approach is not only expensive but also needs human intervention and cannot be considered as an easy automated solution.

Silvestro et al. [18] in 2020, used avatars to allow the user to select the size and the shape, but regarding their results the authors mentioned that they are still working on their project to get good results. Ileperuma et al. [19] in 2020 used the CNNs combined with augmented reality (AR) to detect human body and create an image for a dressed object. The authors claimed to achieve 99% accuracy for their generated images. Nande et al. [20] in 2021 used generative adversarial networks (GANs) to replace the cloth item the customer already wearing, by the cloth item the customer wants to buy in a new image. The authors confirmed that they achieved structural similarity index measure (SSIM) matrices of 0.8.

Chen et al. [21] in 2021 also used M5 transformer to build the VFR and they compared their results to other state-of-the-art DL approaches in literature and they achieved very good results. Lyu et al. [22] in the same year, the authors proposed a High-Resolution Virtual Try-On network (HR-VTON) model to synthesize virtual fitting images, which consists of three sub-modules, namely, a clothing matching module, a try-on module and a refine module. They tested their proposed model on Zalando datset and they confirmed that they achieved accuracy rate of 81%. In 2021 Hyder et al. [23] studied the capabilities of Microsoft Kinect sensor and the role of augmented reality in simulating the surrounding world. The authors confirmed that 85% of the volunteers experienced their system and recommended it as a good 3D learning system. Singh et al. [24] in 2021 also, used Generative Adversarial Networks (GANs) to create a VFR. The authors implemented three models which are: the semantic generation module, the clothes wrapping module, and the content fusion module, and they tested their approach using the Zolando dataset also, and they confirmed that their approach achieved very robust results.

In fact, Lye et al. and Singh et al proved to obtain the best results in literature, but on limited dataset with very narrow variety of images as shown in Section IV(A)(1). Chandani et al. [25] in 2022 tried many methods to build a VFR like AR, ANNs, and CNNs. The authors claimed that the users can move freely infront of camera, and change the colors of the outfit, but they mentioned nothing numeric about the mean error or the accuracy rate. Malathi et al. [26] in 2022 also used DL to reconstruct the 3D graphical perception of the user from his/her 2D image to create cost efficient VFR rather than using Kinect camera. The authors confirmed that their approach achieved very good results, with many useful features that maximize user experience satisfaction.

Prabhakar [27] et al. in the same year, used Alpha Channel Masking to mask the user's shirt to gain an area of interest. The authors only worked on three different colours of the same style (Shirt). Mohamed et al. [28] in 2023 used the 3D reconstruction of human body to build VFR. The authors used Kinect combined with other more expensive body scanners, but they found out that their approach is not only expensive but also needs human intervention and cannot be considered as an easy automated solution. The authors used neural networks to reconstruct the 3D models of the human body and the cloth item and put them together in one photo. The authors also used CNNs combined with ResNet-50 to identify the type of texture to make the reconstructed image more accurate and closer to reality, but in their results, they showed a bit large mean error.

Omkar et al. [29] in 2023 have combined AR with DL to build a VFR. The authors confirmed that their approach that renders the image cloth on the customer's image improved the customer's experience, but they also mentioned nothing factual about their results. Yang et al. [30] in 2023 considered the body mass (BM) of the user and they confirmed that their approach improved the virtual fitting results to a great extent.

## III. METHOD

This work is designed to be deployed on a mobile application in which the process would be much easier for the consumer to do and much cheaper for the clothing shops. The proposed system mainly consists of three main phases: image preprocessing, the GANs part consists of a conditional generator to generate a new segmentation map with the person wearing the new cloth item, and the image generator to generate the final image. The generic flow of the system is as follows: The system asks the user to enter two input images which are, the person image in any pose, and an image for the cloth item he/she wants to buy, to generate a new image of the person wearing the new cloth item. Following is the explanation of every module. Fig. 1 demonstrates the complete architecture of the proposed system.

Fig. 1. The complete architecture of the proposed system.

## A. Images Preprocessing

The input image preprocessing module is composed of four different components: OpenPose, DensePose, segmentation, and cloth mask (see Fig. 1). As mentioned before the user should provide the system by two images, which are his/her image, and the cloth item he/she wants to buy. Regarding the cloth image, it should be segmented using a masking model to be able to recognize the cloth and separate it from the background. Regarding the person's image, some steps should be held such as: background removal, segmentation of different human body parts, definition of OpenPose key points to wrap the cloth over the person in any position not only from the front, 3D human image construction using dense-pose to enhance the final wrapped cloth on the person pose, and finally the combination of the open-pose and the dense-pose to create an agnostic image that would make the system able to focus on the key parts of the human body that should place the new t-shirt on. The outputs of the preprocessing module would be the dense-pose, clothing mask, and segmented agnostic image, which will be the inputs to the conditional generator that generates the new image. Following is a detailed explanation of the preprocessing steps.

*1) Human body segmentation:* Human parts segmentation is the most important preprocessing step for the generator, as the following modules of the proposed system depend on its results. There are 20 parts in the human body image that should be segmented and classified accurately to identify the clothes' location and these classes are: Background, hat, hair, gloves, sunglasses, upper clothes, dress, coat, socks, pants, torso skin, scarf, skirt, face, eft arm, right arm, left leg, right leg, left shoe, and right shoe (see Fig. 2).



Fig. 2. Typical human body segmentation sample.

The Crowd Instance-level Human Parsing (CIHP) model [31] was used to segment the parts of the human body. This model relies on the concept of part grouping network (PGN) [32] where the input image is scaled into six scales which are: 0.5, 0.75, 1, 1.25, 1.5, 1.75 of the original width and height. The input to this model should include the human image its invert. In other words, the input to this model is a stack of image pairs; the image and its invert [image, image invert]. For each input the average of scales is calculated for every input image. This pipeline consists of six models makes the inference time very large (140 sec), but it gives very accurate results, however, this huge inference time is not suitable for this use-case, so we thought how to decrease it. To solve the problem of large inference time we tried other models for segmentation such as self-correction human parsing [33], and Deep-Lab V3 [34]. After utilizing and evaluating the three models Deep-Lab V3 model were selected for this work (see Section V).

*2) Agnostic images generation:* Agnostic images are created to eliminate all the old clothes from the input image and the segmented image. First, a color scale conversion from RGB to gray is done for both input and segmentation images. Then, the background color is converted to black.

This can be done using the open pose key points and the segmentation colors map which help us to define where the upper parts are to eliminate them, and we used the key points to mask the hands only because the segmentation gives the information of where the left and right arms are. Optimization was done using the high-quality segmentation and open-pose key points. Fig. 3 represents the agnostic image generation steps.



Fig. 3.    Agnostic images generation.

After many experiments, we found that eliminating the background gives us more accurate results because the dataset distribution which the model trained on with a white background, The background class is one of the segmentation classes, so we exploit that to eliminate the background by looping on the image and change the color of it to white using open cv (see Fig. 4).



Fig. 4.    Background removal example.

*3) Cloth masking:* Cloth masking is one of the main modules of the proposed system. To build a model that can generate the mask of the clothes; two approaches were suggested in this work, in the first approach, three models were evaluated; AIM, Timi-Net, and P3M model to detect the main object in the image and remove other objects, so, they can be used to generate the mask of the clothes if the clothes considered as the main object in the image. The second approach is cloth segmentation which is used to generate cloth mask in two steps instead of one step first by performing segmentation and then binarizing the result. After

evaluation and analysis of the suggested models, the cloth-segmentation was found to be the champion model which is a robust model against the different colors, and backgrounds [35 - 37]. Fig. 5 shows a sample result for cloth masking.



Fig. 5.    Cloth masking sample result.

*4)Using dense pose:* The fourth and final sub-module in the preprocessing module is the DensePose, in which (R_101_FPN_DL_s1x) [38] was used and succeeded in achieving the best result compared to the baseline model used by Lyu et al [22] (see Section V). Fig. 6 shows the difference between the DensePose proposed by this work, and the DensePose of the work in literature.



Fig. 6.    Enhanced DensePose (right) vs. original or literature's DensePose (middle).

### B. Output Image Generation using Conditional GANs

The second phase of the system is to generate a new image for the customer wearing the new cloth item, GANs was chosen and deployed because of its ability to generate new images with high accuracy. The GANs use a conditional generator where a wrapped cloth and a segmentation map of the person wearing the new cloth are the targets. There are two encoders with five residual blocks; cloth encoder and segmentation encoder, and one decoder that take these targets as conditions to get the most accurate cloth image appearance flow and segmentation map. Then the next fusion block in the encoder takes the previous output after passing to a 3x3 convolutional layer to predict the appearance flow-map from the flow pathway and segmentation features from the seg pathway. Information is exchanged through these two pathways to get the most accurate segmentation map and appearance flow. The last block of the decoder is the conditional aligning that removes the overlapping regions from the cloth mask (convert the straight mask to be fitted on the person image) and handle occlusions (remove hands or any parts from person image in front of the cloth). The following hyperparameters were used: pixel wise cross entropy loss function, Perceptual and L1 losses for best wrapping of the clothes, in addition to least-squared GAN loss. Also, multi-scale discriminators were applied for the conditional adversarial loss calculation.

## C. Image Generator

This phase includes a series of residual blocks, up sampling layers, Spade normalization layers based on the segmentation map parameters. It takes the image agnostic, new segmentation map, dense pose and wrapped close resized and concatenated to the activation. To evaluate the generated image a discriminator rejection was deployed to filter segmentation maps with low quality. It is based on data distribution and implicit distribution from the generator. If the output image is of very low quality, it would be rejected.

## IV. EXPERIMENTAL WORK

We have tested some architectures in literature, and it is found to be working perfectly on females as the authors selected only the females' images from Zolando Dataset, however, accordingly, their approach didn't work with any image outside their dataset as shown in Fig. 7. This fact tells that there is an overfitting in the results found in literature. To avoid the overfitting problem in this work, new images from multiple websites were scrapped and trained beside the dataset found in literature. The following section describes both the dataset in literature and the dataset collected for this work.



Fig. 7. Results obtained from their dataset (top) vs. results obtained using image outside the dataset (down).

## A. Datasets

There are two different datasets used in this work for training and validation, one of them is open-source data used in literature and the other one is collected for this work. The following sub-sections describe the datasets used in this study.

*1) Zolando dataset:* Zolando is a high-resolution dataset that is used for virtual try-on tasks. It consists of the frontal view woman ad top clothing image pairs. It is split into about 11K training pairs and 2K testing pairs. This dataset was used by lye et al. [22], who obtained the best results in literature, and for that reason we used their approach as a baseline model to compare the proposed approach with.

*2) New dataset collected in this work:* About 5K images of males' images from Zara, Farfetch, and Zalando websites were scrapped to avoid the overfitting issue. The collected dataset is publicly available on Kaggle [39]. Some challenges in the collected dataset are tattoos, birth marks, caps, dark sunglasses, beard, head cuts, Black and Asian people, and complicated backgrounds (see Fig. 8). The dataset collected was meant to cover a diverse dataset distribution, and balance between males and females' images as follows: images from Zolando dataset (females), images from Zolando (males),

images from Zara (males), images from Farech (males), and images from Farech (females).



Fig. 8. Some challenges in the new scrapped dataset collected in this work.

## B. Validation and Verification

The system was validated at the submodules level starting from the cloth masking, going through DensePose, OpenPose, segmentation, closing mask, the application itself, and reaching the deployment of the app on AWS server. The integrated submodules were also tested using many test cases and the proposed approach perfectly outperformed the state-of-the-art approaches. Based on the new implementation of the segmentation mask, the response time was reduced from 4 minutes to 78 seconds. The model was deployed on AWS successfully and the application received requests and responded concurrently. Fig. 9 shows the difference between our output (right) and the state-of-the-art result (middle) for the same image. Fig. 10 shows the effect of clothes masking and DensePose, in which it can be noticed that there is a great enhancement of the clothes' alignment.



Fig. 9. The difference between the proposed work output (right) and the state-of-the-art result (middle) for the same images.



Fig. 10. The state-of-the-art results (middle) vs. the proposed approach results (right).

## V. RESULTS AND DISCUSSION

This section demonstrates the results achieved by this proposed work and discusses them. The following subsections show every submodule, the results, and the analysis of these results.

## A. Segmentation

CHIP, Self-correction human parsing, and Deep-Lab V3 segmentation models were used and compared in this work to select the best as mentioned before. The comparison was done

in terms of segmentation results and execution time. The following sub-section describes every model and the findings of using it.

*1) Self-correction human parsing segmentation:* To increase the dependability of both the learned models and true labels, the self-correction model uses a learning scheduler to infer more trustworthy pseudo-masks by repeatedly aggregating the learned model with the former ideal one in an online learning. Using an annotation initialization model as the first step in the learning process, and then adjusting the labels according to the data can improve the model's performance. Through cycles of self-correction learning, both the models and labels improve in accuracy and strength. This model has a good inference time (2 sec), but its pre-trained weights were on a different dataset which missing the torso skin class (the neck class). The pre-trained weights were on the Look into person (LIP) dataset [40] which contains the same classes, but the only difference is the torso-skin class is replaced with a jumpsuit, so the results were without the torso-skin class. To address the problem of the missing torso class, the model was trained on another dataset including all classes, but the results were worse than the CHIP model, and it required more computational resources for training and testing (see Fig. 11).



Fig. 11. Sample of self-correction segmentation results.

*2) Deep-Lab V3 model segmentation:* DeepLab V3 utilizes a novel architecture that combines multi-scale features, dilated convolutions, and conditional random fields to produce high-quality segmentation masks. It is highly effective at handling complex scenes, such as those found in urban environments, and achieves superior performance compared to other state-of-the-art segmentation models. The model was trained to segment all parts within much less time. DeepLab V3 model worked in 0.23 seconds only to produce segmented image. Although the mean Intersection Over Union (IOU) is less than mean IOU of CHIP model, we have chosen DeepLab V3 over CHIP model due to the great difference in execution time, meanwhile, the mean IOU difference between both models is not so large to consider

over the time. Table I demonstrates the results and execution time of the three models. Fig. 12 shows the segmentation results of CHIP and DeepLab V3 models.

TABLE I. SEGMENTATION RESULTS AND EXECUTION TIMES OF THE USED SEGMENTATION MODELS

| Model | Time (seconds) | IOU |
|---|---|---|
| CHIP | 38.41 | 78.21 |
| DeepLab V3 | 0.23 | 64.53 |



Fig. 12. Segmentation results of DeepLab V3 (right) and CHIP (left) models.

*B. Open Pose*

Three different approaches were tested to find the optimal open pose model. Those three models are the Media-pipe model [41], the Detectron model [42], and the CMU-Perceptual-Computing-Lab [43 & 44]. The first model failed to address the occlusion challenge (see Fig. 13), the second model outperformed the results of the first model, but still missing some important points. The third model solved the problems of the previous two models successfully. The following figures show the differences between the results obtained by the three models. Here Fig. 14 shows Detectron open pose results and Fig. 15 shows the CMU-Perceptual open pose detection results.



Fig. 13. Media-pipe pose detection results (occluded arms are missing).



Fig. 14. Detectron open pose detection results (better but some important points are still missing).

Fig. 15. CMU-Perceptual open pose detection results (the best results out of the three models).

*C. Cloth Masking*

To verify the quality of the cloth masking model, three experiments have been held to check the performance of the proposed approaches to get the best cloth masking model. The first experiment was to find the best model of the three models that perform image matting, namely: AIM, Timi-Net, and P3m [45 - 48]. A dataset of 5K images dataset has been scrapped from different fashion websites, then it has been used to compare the results of each one of the three models. The AIM model obtained the worst in terms of average inference time per image, and ability to generate the mask. The abilities of P3M and Timi-Net models to generate the mask were almost the same, but P3M proved to be 1.5 faster than Timi-Net, so, the P3M model was chosen for cloth masking submodule. A manual analysis for all images was then conducted to define the strength and the failure points of the P3M model. As a result of the analysis, P3M showed great results for all cases except in the case of white clothes on white backgrounds, and to solve this problem, image binarization step was added to the segmentation module. The second experiment has been conducted to check the applicability of the second approach to generate the mask using YOLO5 and YOLO7 [49 - 53] for image segmentation, then binarizing the results. The proposed method showed high capability of detecting the existence of the clothes in the images, but at the same time YOlO5 and YOLO7 were not the best choices to perform the cloth segmentation task since they miss some of the images without detecting the existence of the clothes in the image. The third experiment has been held to check the performance of two models which are the U2Net and the cloth-segmentation model. The top 20 images that P3M models failed generating their masks were used in this experiment. U2Net failed with some white clothes with white background, on the other hand, cloth-segmentation model performs well since the model is robust against different colors and backgrounds. As a result, the cloth-segmentation model has been chosen to be integrated with the other models to complete the proposed solution.

*D. Dense Pose*

Different state-of-the-art models in this area were applied, and "Detectron 2" proved to be the best amongst them. Detectron 2 was originally developed by Facebook, where 27 models were trained. All the 27 models were used and evaluated in this work to measure how suitable they are to this problem and select the best accordingly.

*E. Image Generator*

This step is the final step which takes the outputs of all the steps preceding it as an input to generate the final output of the system which should be the image of the user wearing the new cloth item. Freshet Inception Distance (FID) score was used to evaluate the results obtained by this work and as a comparison metric with the work in literature, where FID score is the most modern metric used to measure the distance between real image and equivalent generated one [54 - 57]. Following are the summary of experiments designed for this step:

- The first experiment validates the work of Lye et al [22] as the model addressing the same problem in literature with the best results.

- In the second experiment, we only replaced the DensePose step in literature by the DensePose proposed in this work.

- The third experiment involves changing only the cloth mask step with the proposed new cloth mask method.

- In the fourth experiment we replaced both the DensePose and cloth mask steps with the new proposed methods.

- In the fifth experiment, we only changed the segmentation model, keeping all the literature model as is.

- The sixth experiment includes deploying the proposed model and evaluating it with all its modules integrated together.

Table II demonstrates results of the six experiments, with "new" referring to the proposed models, and "original" referring to the models in literature [22].

TABLE II. COMPARISON BETWEEN THE STATE-OF-THE-ART APPROACHES AND THE PROPOSED APPROACH

| E# | Segmentation | Cloth Mask | Dense Pose | FID |
|---|---|---|---|---|
| 1 | Original | Original | Original | 11.796 |
| 2 | Original | Original | New | 12.243 |
| 3 | Original | New | Original | 11.847 |
| 4 | Original | New | New | **11.743** |
| 5 | New | Original | Original | 13.140 |
| **6** | **New** | **New** | **New** | **11.753** |

As we can see from the above table the proposed approach outperforms the best model in literature, but it can also be noticed that keeping the segmentation of the literature combined with the proposed Cloth mask and DensePose models would produce even better results, i.e., smaller FID.

*F. User Interface*

Finally, to build a complete solution, a user interface has been developed using flutter technology, and the implementation of the application has been divided into two phases. Firstly, creating an initial design by building the main activities, then new features were added using UI/UX. Fig. 16 introduces some of the user interface screens.

Fig. 16. Some of the final user interface screens.

## VI. CONCLUSION AND FUTURE WORK

In conclusion, this work aims at creating a virtual fitting room via a mobile application to make it easier for customers to try on many cloth items without physically dressing them. The best state-of-the-art approach in literature (Lue et al. [22]) has been used in this paper as the baseline architecture to compare the results obtained in this work with. Another reason for choosing them is that they used the Zolando dataset which is publicly available for validation and comparison. But it was found that there is an obvious overfitting because of the unified nature of the dataset. The solution for this included scrapping new data and retraining preprocessing models. A new dataset from different fashion sources was scrapped to collect 5K images to evaluate the proposed approach and compare it to the state-of-the-art approach. All preprocessing subsystems were analyzed and experimented to get the best model. These subsystems included OpenPose (CMU), DensePose (Detectron- R_101_FPN_DL_s1x), cloth mask (cloth segmentation), and segmentation (DeepLab V3). At the end, the proposed approach outperformed the state-of-the-art and succeeded to reduce the FID. The mobile app developed in this work was deployed on AWS service with dockerization technique.

As future work, we suggest adding a new dataset for bottom part and other types of fashion items, in addition to a recommendation system. Other research suggestions may include utilizing newly collected data to retrain the model for improved performance. To reduce inference time semantic segmentation should be considered while training the model, optimizing it for cases with only one instance. Additionally, the use of a teacher-student model (Knowledge Distillation) can be suggested to enhance the performance of all framework models, aiming for better inference times overall. Future work would also include considering children of different ages in testing the proposed architecture.

## ACKNOWLEDGMENT

## REFERENCES

[1] Mătăşel, Alice & Avadanei, Manuela & Talpa, Andreea & Loghin, Carmen. (2023). Virtual Fitting Room and Its Potential in E-Commerce. 10.2478/9788367405133-033.

[2] Mingyu Lu, Suyin Chen, Haotian Lin3, and Yueyi Li, "The Study of Virtual Fitting Room in China", Advances in Economics, Business and Management Research, volume 203 Proceedings of the 2021 3rd International Conference on Economic Management and Cultural Industry (ICEMCI 2021).

[3] Lee, Hanna & Xu, Yingjiao. (2018). 2018 Proceedings Classification of Virtual Fitting Room (VFR) Technology in the Fashion Industry: From the Perspective of Customer Experience.

[4] Pepper, R. M., Freeland-Graves, J. H., Yu, W., Stanforth, P. R., & Xu, B. (2011). Evaluation of a rotary laser body scanner for body volume and fat assessment. Journal of Testing & Evaluation, 39(1), 1–6.

[5] Alfredo, C., and Rodriguez, B. (2016). Virtual fitting rooms (Unpublished master's thesis). Retrieved from oa.upm.es/ 42311/1/TFM_BECERRA_RODRIGUEZ_CARLOS_a.pdf.

[6] Boonbrahm, P., Kaewrat, C., & Boonbrahm, S. (2015). Realistic simulation in virtual fitting room using physical properties of fabrics. Procedia Computer Science, 75, 12–16.

[7] Gao, Y., Brooks, E., & Brooks, A. (2014). The performance of self in the context of shopping in a virtual dressing room system. Proceedings of the International Conference on HCI in Business (pp. 307–315) Heraklion, Greece.

[8] Javornik, A., Rogers, Y., Moutinho, A., & Freeman, R. (2016). Revealing the shopper experience of using a magic mirror' augmented reality make-up application. Proceedings of the 2016 ACM Conference on Designing Interactive Systems (DIS) (pp. 871–882). New York, NY.

[9] Lee, H., & Leonas, K. K. (2018). Customer experiences, the key to survive in an omni-channel environment: Use of virtual technology. Journal of Textile and Apparel, Technology & Management, 10(3), 1–23.

[10] Francisco Pereira, Catarina Silva1, and Mário Alves, "Virtual Fitting Room Augmented Reality Techniques for e-Commerce", M.M. Cruz-Cunha et al. (Eds.): CENTERIS 2011, Part II, CCIS 220, pp. 62–71, 2011. © Springer-Verlag Berlin Heidelberg 2011.

[11] Ioannis Pachoulakis and Kostas Kapetanakis, "AUGMENTED REALITY PLATFORMS FOR VIRTUAL FITTING ROOMS", The International Journal of Multimedia & Its Applications (IJMA) Vol.4, No.4, August 2012.

[12] Dias, Jessica and Chouhan, Divya and Churi, Preshit and Parab, Pranay, Augmented Reality Based Virtual Dressing Room Using Unity3D (April 8, 2022). Available at SSRN: https://ssrn.com/abstract=4111818

[13] N. Z. Hashmi, A. Irtaza, W. Ahmed and N. Nida, "An augmented reality based Virtual dressing room using Haarcascades Classifier," 2020 14th International Conference on Open-Source Systems and Technologies (ICOSST), Lahore, Pakistan, 2020, pp. 1-6, doi: 10.1109/ICOSST51357.2020.9333032.

[14] Jimit Mehta, Priyam Patel, Hiral Katira, and Aarti Sahitya, "Enhancement in Shopping Experience of Clothes Using Augmented Reality", International Journal of Research in Engineering, Science and Management Volume-3, Issue-3, March-2020 www.ijresm.com | ISSN (Online): 2581-5792.

[15] Boonbrahm, Poonpong & Kaewrat, Charlee & Boonbrahm, Salin. (2015). Realistic Simulation in Virtual Fitting Room Using Physical Properties of Fabrics. Procedia Computer Science. 75. 12-16. 10.1016/j.procs.2015.12.189.

[16] França, Ana Carol & Soares, Marcelo. (2018). Review of Virtual Reality Technology: An Ergonomic Approach and Current Challenges. 52-61. 10.1007/978-3-319-60582-1_6.

[17] Rodríguez, Becerra and Carlos Alfredo. "Virtual Fitting Rooms." (2016).

[18] Francesco Sapio, Andrea Marrella, and Tiziana Catarci, "Integrating Body Scanning Solutions into Virtual Dressing Rooms", Preprint of 2018 International Conference on Advanced Visual Interfaces (AVI '18).

[19] Silvestro V. Veneruso, Tiziana Catarci, Lauren S. Ferro, Andrea Marrella, and Massimo Mecella. 2020. V-DOOR: A Real-Time Virtual Dressing Room Application Using Oculus Rift. In Proceedings of the International Conference on Advanced Visual Interfaces (AVI '20). Association for Computing Machinery, New York, NY, USA, Article 99, 1–3. https://doi.org/10.1145/3399715.3399959.

[20] I. C. S. Ileperuma, H. M. Y. V. Gunathilake, K. P. A. P. Dilshan, S. A. D. S. Nishali, A. I. Gamage and Y. H. P. P. Priyadarshana, "An Enhanced Virtual Fitting Room using Deep Neural Networks," 2020 2nd International Conference on Advancements in Computing (ICAC), Malabe, Sri Lanka, 2020, pp. 67-72, doi: 10.1109/ICAC51239.2020.9357160.

[21] Tejas Nande, Khushwant Salve, Gaurav Patange, Swapnil Burde, Shlok Mishra, and Nisarg Gandhewar, "Smartfit - Artificial Intelligence Based Virtual Dressing Room", International Journal of Advanced Research in Science, Communication and Technology (IJARSCT) Volume 4, Issue 3, April 2021.

[22] Jie Chen, Junwen Bu, and Zhiling Huang, "Image-based Virtual Fitting Room", arXiv:2104.04104v1 [cs.CV] 8 Apr 2021.

[23] Hyder, Hasnain & Baloch, G. & Saad, Khawaja & Shaikh, Nehal & Buriro, Abdul Baseer & Ahmed, Junaid. (2021). Particle Physics Simulator for Scientific Education using Augmented Reality. International Journal of Advanced Computer Science and Applications. 12. Doi: 10.14569/IJACSA.2021.0120284.

[24] Qi Lyu, Qiu-Feng Wang, and Kaizhu Huang, "High-Resolution Virtual Try-On Network with Coarse-to-Fine Strategy", Journal of Physics: Conference Series, 1880 (2021) 012009 IOP Publishing doi:10.1088/1742-6596/1880/1/012009.

[25] Rahul Singh, Aman Bindal, Md Azad Khan, Rakshith MR, and Divakara N, "Virtual Clothing Try-on Using Generative Adversarial Networks", IJARIIE-ISSN(O)-2395-4396, Vol-7 Issue-3 2021.

[26] Chandani Lachake, Pooja Badekar, Pravin Shinde, Snehal Kandekar, and Aditi Bharti, "Vastrani (Virtual-Garment-Try-On)", International Journal of Research Publication and Reviews, Vol 3, no 11, pp 2150-2154 November 2022.

[27] M Malathi, R Induja, S Mubarak, and T Lalitha, "Cost Efficient Virtual Trial Rooms", Journal of Physics: Conference Series, 2325 (2022) 012006 IOP Publishing doi:10.1088/1742-6596/2325/1/012006.

[28] T S Prabhakar, N M Shreyas, Akshay Raghu, Chethan B R, and Impana G Shetty, "A Novel Approach of Virtual Visualization of Cloth Fitting", International Journal of Engineering Research in Computer Science and Engineering (IJERCSE) Vol 9, Issue 10, October 2022.

[29] Yasmin Mohamed, Salma Osama, Rogina Michelle, Youssef Karam, Mennat Allah Hassan, and Diaa Salama, "Fit Moi: Online Virtual Fitting Room with Texture Identification and Recommendation System", Journal of Computing and Communication Vol.2, No.2, PP. 19-30, 2023.

[30] Omkar S. Jadhav, Kiran S. Ambatkar, Prathamesh D. Kulkarni, Anuja S. Modhave, Supriya S. Gadekar, "VIRTUAL TRIAL ROOM USING AR AND AI", Journal of Emerging Technologies and Innovative Research (JETIR), March 2023, Volume 10, Issue 3

[31] Yang, S., Xiong, G., Mao, H., & Ma, M. (2023). "Virtual Fitting Room Effect: Moderating Role of Body Mass Index". Journal of Marketing Research, 0(0). https://doi.org/10.1177/00222437231154871

[32] Dandan WANG and Tianci Zhang, "Establishment and Optimization of Video Analysis System in Metaverse Environment", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 14, No. 10, 2023. https://thesai.org/Downloads/Volume14No10/Paper_6-Establishment_and_Optimization_of_Video_Analysis_System.pdf

[33] Manfredi G, Gilio G, Baldi V, Youssef H, Erra U. VICO-DR: A Collaborative Virtual Dressing Room for Image Consulting. *Journal of Imaging*. 2023; 9(4):76. https://doi.org/10.3390/jimaging9040076

[34] Jie Yang, Chaoqun Wang, Zhen Li, Junle Wang, and Ruimao Zhang, "Semantic Human Parsing via Scalable Semantic Transfer over Multiple Label Domains", a parsing network for a specific label domain by employarXiv:2304.04140v1 [cs.CV] 9 Apr 2023.

[35] Gong, K.; Wang, X.; Tan, S. Correlating Edge with Parsing for Human Parsing. Electronics 2023, 12, 944. https://doi.org/10.3390/electronics12040944.

[36] Li, Peike & Xu, Yunqiu & Wei, Yunchao & Yang, Yi. (2020). Self-Correction for Human Parsing. IEEE Transactions on Pattern Analysis and Machine Intelligence. PP. 1-1. 10.1109/TPAMI.2020.3048039.

[37] Zafar, Mehwish & Amin, Javeria & Sharif, Muhammad & Anjum, Muhammad & Mallah, Ghulam & Kadry, Seifedine. (2023). DeepLabv3+-Based Segmentation and Best Features Selection Using Slime Mould Algorithm for Multi-Class Skin Lesion Classification. Mathematics. 11. 364. 10.3390/math11020364.

[38] Tatli, Umut & Budak, Cafer. (2023). Biomedical Image Segmentation with Modified U-Net. Traitement du Signal. 40. 523-531. 10.18280/ts.400211.

[39] Zhang, Ruifei & Liu, Sishuo & Yu, Yizhou & Li, Guanbin. (2023). Self-Supervised Correction Learning for Semi-Supervised Biomedical Image Segmentation. 10.48550/arXiv.2301.04866.

[40] https://github.com/facebookresearch/detectron2/blob/main/projects/DensePose/doc/DENSEPOSE_IUV.md#References

[41] https://www.kaggle.com/datasets/girgismicheal/viton-dataset

[42] https://paperswithcode.com/dataset/lip

[43] Kim, Jong-Wook & Choi, Jin-Young & Ha, Eun-Ju & Choi, Jae-Ho. (2023). "Human Pose Estimation Using MediaPipe Pose and Optimization Method Based on a Humanoid Model". Applied Sciences. 13. 2700. 10.3390/app13042700.

[44] Sahoo, S.K., Palai, G., Altahan, B.R. et al. "An Optimized Deep Learning Approach for the Prediction of Social Distance Among Individuals in Public Places During Pandemic. New Gener". Comput. 41, 135–154 (2023). https://doi.org/10.1007/s00354-022-00202-1.

[45] Nakano, Nobuyasu & Sakura, Tetsuro & Ueda, Kazuhiro & Omura, Leon & Arata, Kimura & Iino, Yoichi & Fukashiro, Senshi & Yoshioka, Shinsuke. (2019). Evaluation of 3D markerless motion capture accuracy using OpenPose with multiple video cameras. 10.1101/842492.

[46] https://cmu-perceptual-computing-lab.github.io/openpose/web/html/doc/md_doc_00_index.html

[47] Yao, Jingfeng & Wang, Xinggang & Ye, Lang & Liu, Wenyu. (2023). Matte Anything: Interactive Natural Image Matting with Segment Anything Models.

[48] Huang, L.; Liu, X.; Wang, X.; Li, J.; Tan, B. Deep Learning Methods in Image Matting: A Survey. Appl. Sci. 2023, 13, 6512. https://doi.org/10.3390/app13116512.

[49] Jizhizi Li, Jing Zhang, and Dacheng Tao, "Deep Image Matting: A Comprehensive Survey", arXiv:2304.04672v1 [cs.CV] 10 Apr 2023.

[50] Ma, S., Li, J., Zhang, J. et al. Rethinking Portrait Matting with Privacy Preserving. Int J Comput Vis **131**, 2172–2197 (2023). https://doi.org/10.1007/s11263-023-01797-8

[51] Oluwaseyi, Olorunshola & Irhebhude, Martins & Evwiekpaefe, Abraham. (2023). A Comparative Study of YOLOv5 and YOLOv7 Object Detection Algorithms. Journal of Computing and Social Informatics. 2. 1-12. 10.33736/jcsi.5070.2023.

[52] T. Reddy Konala, A. Nammi and D. Sree Tella, "Analysis of Live Video Object Detection using YOLOv5 and YOLOv7," *2023 4th International Conference for Emerging Technology (INCET)*, Belgaum, India, 2023, pp. 1-6, doi: 10.1109/INCET57972.2023.10169926.

[53] Juan Terven and Diana Margarita Córdova-Esparza, "A COMPREHENSIVE REVIEW OF YOLO: FROM YOLOV1 AND BEYOND", arXiv:2304.00501v4 [cs.CV] 7 Aug 2023.

[54] Yu, Yu & Zhang, Weibin & Deng, Yun. (2021). Frechet Inception Distance (FID) for Evaluating GANs.

[55] Made Raharja, Surya Mahadi, and Nugraha Priya Utama, "Indonesian Text-to-Image Synthesis with Sentence-BERT and FastGAN", arXiv:2303.14517v1 [cs.CV] 25 Mar 2023.

[56] Tuomas Kynkäänniemi, Tero Karras, Miika Aittala, Timo Aila, and Jaakko Lehtinen, "The Role of ImageNet Classes in Fréchet Inception Distance", arXiv:2203.06026v3 [cs.CV] 14 Feb 2023.

[57] J. Lee and M. Lee, "FIDGAN: A Generative Adversarial Network with An Inception Distance," *2023 International Conference on Artificial Intelligence in Information and Communication (ICAIIC)*, Bali, Indonesia, 2023, pp. 397-400, doi: 10.1109/ICAIIC57133.2023.10066964.

# Observational Quantitative Study of Healthy Lifestyles and Nutritional Status in Firefighters of the fifth Command of Callao, Ventanilla 2023

Genrry Perez-Olivos[1], Exilda Garcia-Carhuapoma[2], Ethel Gurreonero-Seguro[3],
Julio Méndez-Nina[4], Sebastian Ramos-Cosi[5], Alicia Alva Mantari[6]
Programa de Estudios de Enfermería, Universidad de Ciencias y Humanidades, Peru[1, 2, 3, 4]
Image Processing Research Laboratory (INTI-Lab), Universidad de Ciencias y Humanidades, Peru[5, 6]

*Abstract*—**Given the high concern for human health, the aim is to determine the relationship between healthy lifestyles and nutritional status among firefighters of the VCD Callao Ventanilla 2023. This study was conducted in four volunteer fire companies, namely B-75, B-184, B-207, B-232, located in the districts of Ventanilla and Mi Perú. The population consists of 291 personnel, with a sample of 168 participants. It was observed that 58.9% (99) of the participants are under 36 years old, 29.8% (50) are between 36 and 45 years old, and 11.3% (19) are 46 years or older. In terms of gender, 62.5% (105) are males. Regarding the duration of their firefighting service, 70.2% (118) have a maximum of 10 years of seniority. On the other hand, 57.7% (97) of the participants have an unhealthy lifestyle, 40.5% (68) have a healthy lifestyle, and 1.8% (3) have a very healthy lifestyle. Regarding the nutritional status of the firefighters in this study, it was found that 53.3% (89) are overweight, 26.8% are considered normal weight, 19.6% (33) are obese, and 0.6% (1) are underweight. Concerning lifestyles, the study revealed that 57.7% of the participants have an unhealthy lifestyle, 40.5% have a healthy lifestyle, and 1.8% (3) have a very healthy lifestyle. It is worth mentioning that according to Rodríguez C's study, 95.2% of volunteers belonging to the B107 Fire Company lead a healthy lifestyle, while 4.8% do not. Statistically, we can assert that there is no significant relationship between the variable of healthy lifestyles and nutritional status. However, it is observed that there is a direct relationship between nutritional status and age. Likewise, it can be affirmed that more than at least 72.9% of the studied population is overweight, either with overweight or obesity.**

*Keywords—BMI; firemen; lifestyles; excess weight*

## I. INTRODUCTION

A volunteer firefighter, affiliated with an organization focused on fire prevention and control, responds to fires, vehicular accidents, and medical emergencies without charge. Their work involves demanding physical tasks in hostile emergency scenarios [1].

These emergency situations trigger immediate physiological responses in firefighters, including increased heart rate, hyperventilation, heightened oxygen consumption, and sweating. According to a 2019 report from the National Institute for Occupational Safety and Health (NIOSH), cardiac arrest is the leading cause of firefighter fatalities, accounting for 40% of 308 deaths over a decade. Many of these individuals also had coronary artery disease, as per death certificates. These findings highlight a significant percentage of firefighter deaths attributable to cardiovascular diseases [2].

Analyzing past studies on the nutritional and physical activity status of students in various professions reveals a considerable number who do not maintain a lifestyle conducive to their well-being. Predisposing factors to health problems include rapid adaptation to constant changes, lack of physical activity, dietary issues, and overall unhealthy habits [3].

Health is a crucial facet of human life, profoundly impacting one's quality of life, with various factors such as physical activity, diet, and mental well-being playing pivotal roles. The World Health Organization (WHO) noted that over 28% of adults worldwide failed to meet their recommended weekly physical activity levels, equating to at least 150 minutes or 75 minutes of intense activity. Disparities in physical inactivity exist between high- and low-income countries. High-income countries report 26% of men and 35% of women as insufficiently active, compared to 12% of men and 24% of women in low-income countries [4].

The global prevalence of physical inactivity raises significant concerns due to its strong correlation with non-communicable diseases (NCDs), contributing to up to 9% of premature deaths worldwide. Additionally, 6% to 10% of NCDs, such as diabetes, coronary heart disease, colon, and breast cancer, can be attributed to physical inactivity [5].

Obesity and overweight are a growing concern, classified as the 21st century's pandemic by the WHO. Approximately 52% of adults worldwide grapple with these issues [6]. In recent years, global cases of overweight and obesity have tripled, signaling a worrying trend linked to chronic diseases, including Diabetes Mellitus II, cardiovascular disorders, and musculoskeletal problems [6-12].

In 2021, WHO defined overweight and obesity as an excess accumulation of fat detrimental to health due to an energy imbalance between calorie intake and expenditure. In Latin America and the Caribbean, 62.5% of adults, 64.1% of men, and 60.9% of women are overweight or obese. When examining only obesity, it affects 28% of adults, with 26% being men and 31% women, making it the region with the highest prevalence according to the WHO [13].

Peru faces a similar challenge, with 69.9% of adults and 42.4% of young people experiencing overweight or obesity. Poor dietary habits contribute to these figures, with 29% of the population consuming junk food weekly. Fried foods account for 87.1% of weekly consumption [14].

A study by Ramírez J, et al. [15], revealed that 33.6% of Peruvians had abdominal obesity, notably higher in women at 51.2%. Geographic location also influences obesity rates, particularly in cities below 1000 meters above sea level.

Malnutrition is another concern, impacting weight gain and associated with various factors, including age and energy imbalances. A study found 15.3% prevalence of diabetes mellitus II, 43.4% for hypertension, and 17.4% for osteoarthritis [16].

Sleep deprivation, prevalent among firefighters due to demanding schedules, poses health risks. Short sleep patterns, less than six hours, increase the risk of stroke, cancer, coronary heart disease, diabetes, anxiety, depression, and workplace accidents [17-18].

Unhealthy eating habits and sedentary lifestyles are common among firefighters, contributing to obesity. Mandatory physical activity training does not always suffice, further promoting sedentary behavior and weight issues [19-21].

Lifestyle encompasses a person's daily activities, customs, housing, habits, culture, environmental interactions, and interpersonal relationships, all of which are entirely modifiable [22].

Lifestyles can either promote health or pose risks. Healthy lifestyles involve behaviors that enhance well-being, including maintaining a balanced diet, regular physical activity, and abstaining from harmful substances like tobacco and alcohol [23].

Furthermore, lifestyle influences one's physical condition, recreation, and leisure management, which directly impact mental health. Additionally, the consumption of toxic substances can significantly affect overall well-being [24].

Nutritional status refers to the equilibrium between nutrient intake and the body's requirements, determined by the availability of diverse foods and nutritional knowledge [25].

Malnutrition, as defined by WHO, encompasses deficiencies, excesses, and energy imbalances in nutrient intake. It is categorized into three groups: wasting, stunted growth, vitamin and mineral deficiencies, and underweight, with a higher incidence among children [26].

Overweight and obesity, identified as excess fat accumulation detrimental to health, result from an energy imbalance related to calorie intake and expenditure. Nutritional status in adults is typically assessed using the body mass index (BMI), calculated as weight in kilograms divided by height in meters squared (Kg/m²). Classification based on BMI is as follows: underweight (BMI < 18.5), normal weight (BMI 18.5-24.99), overweight (BMI 25.00 – 29.99), and obesity, further classified as class I (BMI: 30.00-34.99), class II (BMI: 35.00-39.99), and class III (BMI ≥ 40.00) [27-29].

In the 21st century, global obesity rates have tripled, labeled as a pandemic by the WHO, a concern not exclusive to Peru. This issue correlates with non-communicable diseases like Diabetes Mellitus II, cardiovascular ailments, and musculoskeletal disorders, affecting a significant portion of the Peruvian population, including firefighters. Regrettably, there's a dearth of studies on Peruvian firefighters' nutritional status, often hampered by small sample sizes that fail to depict the true extent of the issue.

Examining the lifestyles and nutritional well-being of Callao firefighters assumes paramount significance due to evident excess weight and associated unhealthy habits. Understanding the interplay of these variables is essential for implementing effective preventive measures. Firefighters must maintain optimal health for enhanced efficiency and to mitigate health risks. Dietary behavior emerges as a pivotal, modifiable risk factor for occupational diseases. Numerous studies have investigated dietary interventions among firefighters, seeking effective methods to improve their dietary habits.

Addressing this behavioral issue holds great societal potential for reducing associated problems. Consequently, this study aims to establish the relationship between healthy lifestyles and nutritional status among firefighters at VCD Callao Ventanilla in 2023.

## II. RELATED JOBS

Arrieta J, Solís I. [26], Costa Rica in 2020 carried out a study entitled "Eating habits, nutritional status and cardiovascular risk in firefighters from 20 to 59 years of age of the XII battalion, Costa Rica, 2020." The objective of this study was to relate eating habits and nutritional status according to body mass index with cardiovascular risk using the Framingham Heart Study formula in firefighters aged 20 to 59 years from Battalion XII, Costa Rica, in 2020. Where a structured interview was conducted with the firefighters participating in the research, with the data obtained, a statistical analysis was carried out with the models that best predict cardiovascular risk among the established variables. In the results obtained, it was evident that only 31% of the firefighters have a normal body mass index, the remaining percentage is overweight or obese type 1. Most of the sample had a low cardiovascular risk according to abdominal circumference and according to the Framingham Heart Study calculator. They conclude that: There is an association between the consumption of alcohol, semi-skimmed dairy products, saturated fats and refined cereals, and cardiovascular risk, at different frequencies of consumption; In addition, type I obesity and the use of the frying cooking method are also associated with this risk.

Echeverria M., [30]. In Ecuador in 2021 a study entitled "Nutritional status and eating habits of the fire department personnel of the Otavalo Canton 2021" was carried out with the objective of evaluating the nutritional status and eating habits of the personnel of the Otavalo Canton Fire Department 2021, a descriptive cross-sectional research was carried out and 31 people were worked with, An online survey was applied to those who were given two sociodemographic variables and eating habits, and anthropometric measurements

such as height and weight were taken to obtain BMI. The main results highlight that the majority of the staff are male, between the ages of 20 and 49 years. 55% of the staff is overweight and 16% obese, 29% are in normal nutritional status, eating habits are inadequate.

Camargo F, Jardim T, Rocha L, Zandonade E, Nívea K. [31]. In Brazil 2020 they carried out an article on "Prevalence of obesity in Brazilian firefighters and the association of central obesity with personal, occupational and cardiovascular risk factors: a cross-sectional study" where the studied population was 1018 firefighters, leaving 892 firefighters who met the inclusion criteria. The main results were that 48.65% of the firefighters were overweight and 10.99% were obese. In terms of body fat percentage, 26.23% of participants were considered obese, while 18.61% of firefighters were considered centrally obese or at risk in the waist circumference measurement.

Chuquipoma J. [32] Lima – Peru in 2019 presented a study on "relationship between previous knowledge in nutrition and nutritional status in firefighters of the company "Salvadora Lima n° 10", 2018" with the aim of determining the relationship between the level of previous knowledge in nutrition and the nutritional status in firefighters of said company. In a descriptive, correlational, cross-sectional study, the results correspond to a total of 50 firefighters evaluated, of both sexes whose ages ranged from 20 to 59 years. This obtained the following results: 72% of the firefighters have a low level of knowledge and 28% have a fair level of knowledge. Regarding the nutritional status of the firefighters, it was observed that 34% had a nutritional status within normal parameters, 44% were overweight and 22% were obese.

Rodríguez C. [33] Chimbote – Peru in 2018 published a paper entitled "Lifestyles and Biosociocultural Factors of the Volunteer Workers of the B-107 Nuevo Chimbote Fire Company, 2017" with the aim of determining the relationship between the lifestyle and biosociocultural factors of the volunteer workers of the fire company of said company, The sample consisted of 42 members whose results were: There is no statistically significant relationship between lifestyle and biosociocultural factors: age, sex. Whether there is a statistically significant relationship between lifestyle: level of education and income. Religion, marital status and occupation could not be linked.

This information alerts us to the problem and the risk that firefighters have to suffer from non-communicable diseases such as cardiovascular diseases, diabetes mellitus II, among others related to poor nutrition and poor lifestyle. That is why control measures must be taken to prevent these diseases through health promotion.

## III. MATERIALS AND METHODS

### A. Research Approach and Design

The present study is quantitative, correlational and cross-sectional in terms of methodological design. Descriptive because it measures, evaluates and collects data on variables, both lifestyle and nutritional status; it is also correlational because it investigates the relationship and appendages

between lifestyles and nutritional status; analytical because it allows the establishment of an association relationship between variables. And finally, it is a cross-sectional research because the data obtained were collected in a specific space and time [22].

### B. Population, Sample and Sampling

In the present study, we worked with a finite population, which was made up of the active firefighters of the Callao VCD, belonging to the third that is made up of four fire companies B-75, B-184, B-207 and B-232. The sum of its personnel makes a population of 291 firefighters. They were selected according to the following criteria.

*1) Inclusion criteria*

- The participant belongs to any of the four companies of the third Brigade of the VCDC.
- The participating personnel are active personnel according to the RIF between 18 and 65 years of age.
- The participant must accept his/her participation in this study voluntarily, and must sign the informed consent form.

*2) Exclusion criteria*

- All candidates who do not meet 100% of the inclusion requirements were excluded.

*3) Sample size:* To determine this sample, the statistical program "EPIDAT" was calculated, whose data to calculate will be a population of 295 effective, with a confidence level of 95.0%, the sample size of 168 participants of this study was obtained.

*4) Sample selection:* The sampling was non-probabilistic due to convenience, ease of access and time availability.

### C. Study Variable(s)

In the following research, lifestyle is the main variable, it is a qualitative variable and its measurement scale is ordinal. On the other hand, the nature of the nutritional status variable will be obtained through weight and height, both of which are quantitative variables.

*1) Conceptual definition of lifestyles:* Lifestyle is the general way of life of each person, the way in which they conduct their day-to-day activities, this can be expressed in behaviors, basically in customs, and it is also configured by housing, habits, culture, relationships with the environment and relationships between individuals.

*2) Operational definition: Lifestyles* is the generic form that is equivalent to the form, mode and manner of living, epidemiologically speaking, lifestyle is the habit of life or the way of life, which would be a set of behaviors that people choose, which can be healthy or harmful to health. In this sense, lifestyles are behaviors that improve or increase health risk [19].

*3) Conceptual definition of nutritional status:* Nutritional status is the balance that exists between the nutritional contribution of what we ingest and the nutritional demands

that our body desires, which determines the quality of nutrients and their use. Among the factors that determine nutrition, it depends on the variety of foods available and access to them. On the other hand, people's level of knowledge in nutrition also influences nutritional status [20].

*4) Operational definition:* Nutritional status was determined using the Body Mass Index (BMI), which was calculated by dividing weight in kilograms by height in meters squared. The following results will be considered: Underweight if the BMI <18.5, normal weight BMI 18.5-24.99, Overweight if the BMI is 25.00 – 29.99 and obese if the BMI is ≥ 30. Obesity will be classified as: class I obesity (BMI: 30.00-34.99), class II obesity (BMI: 35.00-39.99) and class III obesity (BMI ≥ 40.00). For this study, excess weight has been considered overweight and obese [28].

## D. Measuring Technology and Instrument

*1) Data collection technique:* The technique used during data collection for the lifestyle variable was the survey, which is widely used in quantitative and descriptive studies. Regarding the variable of nutritional status, we used the measurement of anthropometric measurements in the field.

*a) For the Lifestyles variable:* The modified and validated questionnaire of Palomares L [25] was used. This lifestyle questionnaire consists of a total of 48 questions, which is divided into 6 dimensions:

- Fitness, Physical Activity & Sport
- (4 questions).
- Recreation and Leisure Time Management
- (6 questions).
- Use of Alcohol, Tobacco, and Other Drugs
- (6 questions).
- Dream (6 questions).
- Eating Habits (18 questions).
- Self-Care and Medical Care (8 questions).

All dimensions correspond to a Likert scale (Never:0, Sometimes:1, Frequently:2, and Always:3).

The rating that was considered to assess healthy lifestyles is:

- Unhealthy: 0 - 36 points
- Unhealthy: 37 - 72 points
- Healthy: 73 - 108 points
- Very healthy: 109 - 144 points

*b) For the Nutritional Status variable:*

*c) To obtain the body weight,* the measurement was carried out with a rechargeable electronic scale of the SEC brand with a capacity of 180kg, model SEC-180 platform, whose results will give us in kilograms (kg).

*d)* To obtain the carving we used a wooden height meter with a measuring range of up to 3.5 to 230 centimeters. It was placed on a smooth and flat surface, without any unevenness or any foreign object under it, and with the board resting on a flat surface forming a right angle with the floor and wall.

Once weight and height were obtained, BMI was calculated using the following expression. BMI = Weight/Height2. Depending on the result obtained, nutritional status was determined by BMI according to WHO criteria. thinness (BMI < 18.5), weight (BMI: 18.5-24.99), overweight (BMI: 25.00-29.99), obesity class I (BMI: 30.00-34.99), obesity class II (BMI: 35.00-39.99) and obesity class III (BMI ≥ 40.00).

## E. Procedure for Data Collection

*1) Authorization and prior coordination for data collection:* For the development of field work, permission will be requested from the chiefs of the four companies of the third Brigade of the VCDC to allow us to enter their units to carry out data collection.

*2) Application of data collection instrument(s)*

*a) For the Lifestyles variable:* A modified and validated questionnaire from Palomares L [25] will be used to measure lifestyles. The questions will be put on several printed sheets so that the participant can easily fill in according to his/her criteria (ANNEX B).

- Explain to participants about the work to be done and ask them to fill out the informed consent document.
- Briefly explain the six dimensions or parts of the questionnaire.
- Guide participants from the beginning to the end of filling out the form.
- Verify the correct completion of the questions and store it for later analysis.

*b) For the variable Nutritional Status: For* the measurement of weight and height, the Ministry of Health's technical guide for the anthropometric nutritional assessment of adults was used as a study in [28]. The data obtained were recorded in the anthropometric data collection form (ANNEX C).

*3) Weight taking*

- The location and condition of the balance was verified. This should be on a smooth, flat, horizontal surface. No unevenness and no presence of any foreign object under the scale.
- Ask the adult to take off his or her shoes.
- Verify that the scale reads 00 (zero) andes of the weight intake.
- Ask the elderly person to stand in the center of the scale platform, in an upright and relaxed position, facing the front of the scale, with the arms at the sides of the body, with the palms resting on the thighs, the heels slightly apart and the tips of the feet apart forming a "V".

- Read the weight in kilograms and decimals expressed in grams, and then deduct the weight of people's garments.

- Record the weight obtained in kilograms (kg), with a decimal place corresponding to 100 g, in the corresponding format, in clear and legible handwriting, (example: 65.1 kg).

*4) Size taking:* Check the wooden height gauge, the sliding of the moving stop must be smooth and without swaying, the tape measure must be well adhered to the board and its numbering must be clearly observed. Likewise, the stability conditions of the tachymeter must be checked.

- Check the location and condition of the tachymeter. Check that the moving stop slides smoothly, and check the condition of the tape measure for a correct reading.

- Explain to the elderly the procedure for measuring the height, slowly and patiently, and ask for their collaboration.

- Ask you to remove shoes (flip flops, sandals, etc.), excess clothing, and accessories or other objects on your head that interfere with the measurement.

- Indicate and help you to position yourself in the center of the base of the tachymeter, with your back to the board, in an upright position, facing forward, with your arms at your sides, with your palms resting on your thighs, heels together and the balls of your feet slightly apart.

- Make sure your heels, calves, buttocks, shoulders, and back of your head are in contact with the dashboard of the height meter

- Check the "map of Frankfurt". In some cases, it will not be possible due to problems with curvature in the spine, injuries or other problems.

- Place the open palm of your left hand on the chin of the elderly person to be carved, then close it gently and gradually without covering the mouth, in order to ensure the correct position of the head on the height meter.

- Slide the movable stop with your right hand until it makes contact with the top surface of the head, slightly compressing the hair; then slide the movable stop upwards. This procedure (measurement) must be performed three times consecutively, moving the moving stop closer and farther away. Each procedure has a value in meters, centimeters, and millimeters.

- Read the three measurements obtained, obtain the average and record it in centimeters with an approximation of 0.1 cm.

*F. Methods of Statistical Analysis*

Once the data of the nutritional status BMI and the result of the filling of the instrument have been obtained, the results of the questionnaire have been obtained. This information was entered into a Microsoft Excel XP software database version 2016. Frequency tables and statistical software such as IBM SPSS Statistics 27 were used to perform the data analysis.

*G. Ethical Aspects*

To carry out this work, the basic concepts of bioethics were taken into account, such as autonomy, non-maleficence, beneficence and justice, for the protection of the participants' data, as well as informed consent will be applied to give clear and precise information to the participants [33].

*1) Principle of autonomy:* This principle refers to the freedom of decision of the participants, so it must be respected by this research. This principle will be applied in this research, for all participants in which they will be asked to sign the informed consent form and it will be through documents that the participation of the participants will be reflected [34].

*2) Principle of beneficence:* In this principle we refer to not causing harm to others, in which we are subject to avoiding harm and doing good to others. In this study, they will be informed of the importance of the study and the results obtained to improve their health and quality of life [35].

*3) Principle of non-maleficence:* This principle deals with not harming the person and having the obligation to decrease the risk of causing harm. In this paper it will be explained that participation in this research will not cause any harm to your health [36].

*4) Principle of Justice:* The principle of Justice refers to the operational part of research ethics, including non-discrimination of study participants. The participants of this research will be treated equally without any discrimination or favoritism [31], [37].

## IV. RESULTS

Table I shows that 58.9% (99) of the participants were under 36 years of age, 29.8% (50) were between 36 and 45 years of age, and 11.3% (19) were 46 years or older. In terms of sex, 62.5% (105) are male. Regarding marital status, 69.0% (116) were single, 27.4% (46) were married, 2.4 (4) were divorced and 1.2% (2) were widowed. In terms of job category, 40.5% (68) are university professionals, 26.2% (44) are technical professionals, 19.0% (32) are university students, 10.7% (18) are self-employed, 1.8% (3) are housewives and 1.8% (3) are military personnel. In terms of how long they have been firefighters, 70.2% (118) have a maximum of 10 years of service, 22.6% (38) from 11 to 20 years and 7.1% (12) more than 20 years of service. Regarding the hierarchical grade of the participants, 47.0% (79) are sectional, 18.5% (31) have the rank of second lieutenant, 15.5% (26) of lieutenant, 8.9% (15) of captain, 7.1% (12) of lieutenant brigadier and 3.0% (5) of the rank of brigadier. On the other hand, 49.4% (83) have donated blood at some point and 58.3% (49) of this group have last donated blood one year ago.

From Table II we can see that the p-value (Sig.) associated with Spearman's Rho correlation coefficient level is 0.767, which is greater than 0.05 (significance level or alpha value); Therefore, there is no statistically significant reason to say that the variables lifestyle and nutritional status are related.

TABLE I.        SOCIO-DEMOGRAPHIC DATA OF THE FIREFIGHTERS OF THE FIFTH COMMAND OF CALLAO

| Socio-demographic data | n=168 | |
|---|---|---|
| | fi | % |
| **Age** | | |
| 17 - 25 | 43 | 25,6 |
| 26 - 35 | 56 | 33,3 |
| 36 - 45 | 50 | 29,8 |
| 46 - 55 | 12 | 7,1 |
| 56 - 65 | 7 | 4,2 |
| **Sex** | | |
| Female | 63 | 37,5 |
| Male | 105 | 62,5 |
| **Marital status** | | |
| Bachelor | 116 | 69,0 |
| Married | 46 | 27,4 |
| Divorced | 4 | 2,4 |
| Widower | 2 | 1,2 |
| **Job Category** | | |
| Housewife | 3 | 1,8 |
| Undergraduate | 32 | 19,0 |
| Military | 3 | 1,8 |
| Technical Professional | 44 | 26,2 |
| University professional | 68 | 40,5 |
| Self-employed | 18 | 10,7 |
| **Service Time** | | |
| Less than 5 years | 57 | 33,9 |
| Ages 5 to 10 | 61 | 36,3 |
| From 11 to 20 years old | 38 | 22,6 |
| From 21 to 30 years old | 10 | 6,0 |
| More than 30 years | 2 | 1,2 |
| **Hierarchical Grade** | | |
| Brigadier | 5 | 3,0 |
| Lieutenant Brigadier | 12 | 7,1 |
| Captain | 15 | 8,9 |
| Lieutenant | 26 | 15,5 |
| Second Lieutenant | 31 | 18,5 |
| Sectional | 79 | 47,0 |
| **Have you ever donated blood** | | |
| Yes | 83 | 49,4 |
| No | 85 | 50,6 |
| **Last Blood Donation Time (months)** | | |
| Up to 6 months | 29 | 34,5 |
| 7 to 12 months | 20 | 23,8 |
| 13 to 18 months | 9 | 10,7 |
| 19 to 24 months | 10 | 11,9 |
| More than 24 months | 16 | 19,0 |

TABLE II.        RELATIONSHIP OF LIFESTYLES AND NUTRITIONAL STATUS

| | | | Healthy Lifestyles | Nutritional status |
|---|---|---|---|---|
| Spearman's Rho | Healthy Lifestyles | Correlation coefficient | 1,000 | ,023 |
| | | Follow-up (bilateral) | . | ,767 |
| | | N | 168 | 168 |
| | Nutritional status | Correlation coefficient | ,023 | 1,000 |
| | | Follow-up (bilateral) | ,767 | . |
| | | N | 168 | 168 |

TABLE III.        QUALITY OF LIFESTYLES

| Healthy Lifestyles | Frequency | Percentage |
|---|---|---|
| **Unhealthy** | 0 | 0,0 |
| **Unhealthy** | 97 | 57,7 |
| **Healthy** | 68 | 40,5 |
| **Very healthy** | 3 | 1,8 |
| **Total** | 168 | 100,0 |

Table III presents the levels of quality of lifestyles, where we show that 57.7% (97) of the participants have an unhealthy lifestyle, 40.5% (68) a healthy lifestyle and 1.8% (3) a very healthy lifestyle.

TABLE IV.        NUTRITIONAL STATUS

| | Frequency | Percentage |
|---|---|---|
| **Low weight** | 1 | 0,6 |
| **Normal weight** | 45 | 26,8 |
| **Overweight** | 89 | 53,0 |
| **Obesity** | 33 | 19,6 |
| **Total** | 168 | 100,0 |

Table IV shows the nutritional status of the firefighters who participated in this research, from which we found that 53.3% (89) are overweight, 26.8% are considered normal weight, 19.6% (33) are obese and 0.6% (1) are underweight.

TABLE V.        ASSOCIATION BETWEEN LIFESTYLES AND SOCIO-DEMOGRAPHIC DATA

| Socio-demographic data | Lifestyles | |
|---|---|---|
| | Value* | p_valor (Sig.) |
| **Age** | 9,493 | 0,302 |
| **Sex** | 0,662 | 0,718 |
| **Marital status** | 4,493 | 0,610 |
| **Job Category** | 5,123 | 0,883 |
| **Service Time** | 6,360 | 0,607 |
| **Hierarchical Grade** | 7,383 | 0,689 |
| **Have you ever donated blood** | 4,423 | 0,110 |

Value = of the Chi-square statistic; p_valor=probability value used to contrast with the significance level (α=0.05).

Table V presents the association between lifestyles and different socio-demographic data such as age, sex, marital status, and job category, length of service, hierarchical grade and whether you have ever donated blood. We observed that the p-value is higher than the significance level (0.05) for all pairs of variables tested, therefore, we conclude that there is no relationship between the socio-demographic variables and the lifestyle of the participants.

TABLE VI. ASSOCIATION BETWEEN NUTRITIONAL STATUS AND SOCIO-DEMOGRAPHIC DATA

| Socio-demographic data | Nutritional status | |
|---|---|---|
| | Value* | p_valor (Sig.) |
| Age | 24,079 | 0,020* |
| Sex | 6,075 | 0,108 |
| Marital status | 12,969 | 0,164 |
| Job Category | 6,724 | 0,965 |
| Service Time | 18,679 | 0,097 |
| Hierarchical Grade | 12,465 | 0,644 |
| Have you ever donated blood | 6,013 | 0,111 |

Value = of the Chi-square statistic; p_valor=probability value used to contrast with the significance level ($\alpha$=0.05).

Table VI presents the association between the nutritional status and the socio-demographic data of the respondents. We found that the p-value (Sig.) is less than 0.05 (level of significance) only in the case of age and nutritional status; this indicates that there is statistically significant evidence to conclude that there is an association between age and nutritional status.

## V. DISCUSSION

It is observed that the collaborators are found a higher percentage: 58.9% are under 36 years of age and 29.8% range between 36 and 45 years of age, and as for sex, it has a higher prevalence in men that is equivalent to 62.5%, according to marital status it is found that singles stand out with 69.0% and married with 27.4%. With regard to the level of education, 40.5% are university professionals and 26.2% are technical professionals. However, [24] found between ages ranging from 20 to 29 years of age, counting 61.3% and 38.7% including those aged 40 to 64 years. On the other hand, 77.4% of males belong to them and 22.6% of females. Meanwhile, the marital status related to married people is 71% and single people found 16.1% are single people.

In terms of how long they have been firefighters, 70.2% (118) have a maximum of 10 years of service, 22.6% (38) from 11 to 20 years and 7.1% (12) more than 20 years of service. Regarding the hierarchical grade of the participants, 47.0% (79) are sectional, 18.5% (31) have the rank of second lieutenant, 15. Lieutenant br5% (26) of lieutenant, 8.9% (15) of captain, 7.1% (12) of l brigadier and 3.0% (5) of the rank of brigadier. On the other hand, 49.4% (83) have donated blood at some point and 58.3% (49) of this group have last donated blood one year ago.

Table III presents the levels of quality of lifestyles, that is, the way we have been consuming our food is there is an orderly or disordered way to acquire it [26] where we

visualize that 57.7% of the participants have an unhealthy lifestyle, 40.5% a healthy lifestyle and 1.8% (3) a very healthy lifestyle. It is worth mentioning that [27], in his study it is evident that of the 100% of the volunteers who belong to the B107 Fire Company, lead a healthy lifestyle is 95.2%, those who do not have a non-soluble lifestyle is the result of 4.8%. It is evident that there is no similarity in our study, because the results are very different. Table IV shows the nutritional status of the firefighters who participated in this research.

In our study on the nutritional status of firefighters, it is evident that overweight is in first place with 53.3%, followed by normal with 26.8% and obesity is 19.6% (5). There are several studies that are almost similar to our study (see Table IV) as evidenced by the result of [26], that overweight results in 44%, on the other hand, normal weight is evidenced in 34% and firefighters are only observed to have obesity in 22%).

## VI. CONCLUSION

The present study examined the lifestyles and nutritional status of the firefighters from VCD Callao Ventanilla in 2023, highlighting key demographic aspects such as age, gender, marital status, education, and years of service. These data provide valuable context for understanding health behaviors in this population.

A high prevalence of unhealthy lifestyles among firefighters is evident, with 57.7% exhibiting concerning dietary habits. Given that dietary patterns influence occupational health and overall well-being, it is crucial to address these habits through specific interventions to improve their health and performance.

Regarding nutritional status, it is observed that a substantial portion of firefighters is overweight (53.3%), followed by 26.8% with a normal BMI and 19.6% classified as obese. While there is some agreement with previous research, significant variations are identified, emphasizing the need for further investigations to understand the factors contributing to these differences in nutritional status.

For future research, longitudinal studies are recommended to track changes in the lifestyles and nutritional status of firefighters over time, providing insights into the effectiveness of interventions and the evolution of health issues. Additionally, it is essential to investigate specific risk factors contributing to unhealthy lifestyles and overweight among firefighters, allowing for tailored interventions to address the root causes of these problems.

Furthermore, the importance of developing and testing intervention programs focused on promoting healthy eating habits and physical activity among firefighters is highlighted, prioritizing the improvement of their overall health and the reduction of occupational risks. Finally, conducting comparative studies with firefighter populations in different regions will help identify geographical variations in lifestyles and nutritional status, facilitating the implementation of more specific and effective interventions.

## REFERENCES

[1] El peruano. Decreto Supremo Que Aprueba El Reglamento Del Decreto Legislativo N° 1260, Decreto Legislativo Que Fortalece El Cuerpo

General De Bomberos Voluntarios Del Perú Como Parte Del Sistema Nacional De Seguridad Ciudadana Y Regula La Intendencia Nacional De Bomberos Del Perú [ Internet]. 2021[Citado El 19 De octubre De 2022]. Disponible. Https://Busquedas.Elperuano.Pe/Normaslegales/Decreto-Supremo-Que-Aprueba-El-Reglamento-Del-Decreto-Legisl-Decreto-Supremo-N-019-2017-In-1534348-3/

[2] Instituto Nacional Para La Seguridad Y Salud Ocupacional Prevención De Muertes Entre Bomberos Por Ataques Cardíacos Y Otros Episodios Cardiovasculares Agudos. Centros Para El Control Y La Prevención De Enfermedades [ Internet]. 2019 [Actualizado El 2022; Citado 19 De Octubre2022]. Disponible: Https://Www.Cdc.Gov/Spanish/Niosh/Docs/2007-133_Sp/

[3] Rosales Y, Cordovez S, Fernández Y, Álvarez S. Estado nutricional y actividad física en estudiantes universitarios. Una revisión sistemática. Rev. Chil. Nutr. [Internet]. 2023 [Citado el 22 de noviembre del 2023];50(4):445-456. Disponible en: https://www.scopus.com/record/display.uri?eid=2-s2.0-85175187972&origin=resultslist&sort=plf-f&src=s&sid=385c9d59ae086de7f7064913c3310100&sot=b&sdt=b&s=TITLE-ABS-KEY%28estado+AND+nutricional+AND+actividad%29&sl=40&sessionSearchId=385c9d59ae086de7f7064913c3310100

[4] Organización Mundial De Salud. Obesidad Y Sobrepeso [Internet].Ginebra [Actualizado El 9 De Junio Del 2021;Citado El 20 De Octubre Del 2022].Disponible En: Https://Www.Who.Int/Es/News-Room/Fact-Sheets/Detail/Obesity-And-Overweight

[5] Organización De Salud. Actividad Física Sobrepeso [Internet]: Oms.2022[Actualizado 05 octubre Del 2022.; Citado El 21 octubre Del 2022]. Disponible En: Https://Www.Who.Int/Es/News-Room/Fact-Sheets/Detail/Physical-Activity

[6] Organización Mundial De Salud. Enfermedades No Transmisibles. [Internet].Ginebra [Actualizado 16 Setiembre Del 2022 ;Citado El 20 De octubre Del 2022].Disponible En: Https://Www.Who.Int/Es/News-Room/Fact-Sheets/Detail/Noncommunicable-Diseases

[7] Ricardo Leon Ayala, Sebastian Ramos Cosi, and Laberiano Andrade-Arenas, "Design of a Mobile Application to Improve the Lifestyle of Patients with Diabetes," International Journal of Interactive Mobile Technologies, vol. 17, no. 5, pp. 100–116, 2023.

[8] Meyluz Monica Paico Campos, Sebastian Ramos-Cosi, Laberiano Andrade-Arenas, "SAFE Mobile Application: Prevention of Violence Against Women," International Journal of Engineering Trends and Technology, vol. 71, no. 12, pp. 299-307, 2023. Crossref, https://doi.org/10.14445/22315381/IJETT-V71I12P228

[9] Alva Mantari Alicia, Arancibia-Garcia Alexander, Chávez Frías William, Cieza-Terrones Michael, Herrera-Arana Víctor and Ramos-Cosi Sebastian, "Abnormal Pulmonary Sounds Classification Algorithm using Convolutional Networks" International Journal of Advanced Computer Science and Applications(IJACSA), 12(6), 2021. http://dx.doi.org/10.14569/IJACSA.2021.0120645

[10] S. Ramos-Cosi and N. I. Vargas-Cuentas, "Prototype of a system for quail farming with arduino nano platform, DHT11 and LM35 sensors, in Arequipa, Peru," International Journal of Emerging Technology and Advanced Engineering, vol. 11, no. 11, pp. 140–146, Nov. 2021, doi: 10.46338/IJETAE1121_16

[11] Sebastian Ramos-cosi, Lina Cardenás-Pineda, David Llulluy-Nuñez, Alicia Alva-Mantari , "Development of 3D Avatars for Inclusive Metaverse: Impact on Student Identity and Satisfaction using Agile Methodology, VRChat Platform, and Oculus Quest 2 ," International Journal of Engineering Trends and Technology, vol. 71, no. 7, pp. 1-14, 2023. Crossref, https://doi.org/10.14445/22315381/IJETT-V71I7P201

[12] G. L. Bakris et al., "Effect of Finerenone on Chronic Kidney Disease Outcomes in Type 2 Diabetes," New England Journal of Medicine, vol. 383, no. 23, pp. 2219–2229, Dec. 2020, doi: 10.1056/NEJMOA2025845.

[13] Organización Panamericana De Salud Y Organización Mundial De Salud Prevención De La Obesidad. [Internet]. Ops/Oms [Citado El 20 De octubre Del 2022]. Disponible En: https://www.paho.org/es/temas/prevencion-obesidad

[14] Instituto Nacional Del Niño Y Ministerio De Salud. Cerca Del 70%De Adultos Peruanos Padecen De Obesidad Y Sobrepeso [Internet]:

Ins/Minsa.2019[Actualizado 28 De Marzo Del 2019;Citado El 21 Octubre Del2022].Disponible En: Https://Web.Ins.Gob.Pe/Index.Php/Es/Prensa/Noticia/Cerca-Del-70-De-Adultos-Peruanos-Padecen-De-Obesidad-Y-Sobrepeso

[15] Pajuelo J, Torres L, Agüero R, Bernui I. Overweight, obesity and abdominal obesity in the adult population of Peru. An. Fac. med. [Internet]. 2019 Ene [citado 2022 Nov 17] ; 80( 1 ): 21-27. Disponible en: http://www.scielo.org.pe/scielo.php?script=sci_arttext&pid=S1025-55832019000100004&lng=es. http://dx.doi.org/10.15381/anales.v80i1.15863.

[16] Ministerio de salud, El Costo De La Doble Carga De La Malnutrición. [Internet].Ginebra [Julio Del 2022 ;Citado El 20 De Octubre Del 2022].Disponible en Https://Docs.Wfp.Org/Api/Documents/Wfp-0000140902/Download/?_Ga=2.189255020.636623899.1666918640-491367601.1666918640

[17] Bernui I, Delgado D. Factores Asociados Al Estado Y Al Riesgo Nutricional En Adultos Mayores De Establecimientos De Atención Primaria. An. Fac. Med. [Internet]. Prensa 2022 [ Citado El 28 De Octubre Del 2022]; 82(4 ): 261-268. Disponible en: Http://Dx.Doi.Org/10.15381/Anales.V82i4.20799.

[18] Instituto Nacional De Estadística E Informática. El 12,1% De La Población Menor De Cinco Años De Edad Del Apis Sufrió Desnutrición Crónica En El Año 2022[Internet]. Prensa 2022 [ Citado El 28 De octubre Del 2022]. Disponible En: Https://M.Inei.Gob.Pe/Prensa/Noticias/El-121-De-La-Poblacion-Menor-De-Cinco-Anos-De-Edad-Del-Pais-Sufrio-Desnutricion-Cronica-En-El-Ano-2020-12838/#:~:Text=En%20el%20a%C3%b1o%202020%2c%20el,De%20resultados%20de%20los%20programas

[19] Díaz E, Rubio S, López M, Aparicio M. Los Hábitos De Sueño Como Predictores De La Salud Psicológica En Profesionales Sanitarios. An. Psicol. [Internet]. 2020 [Citado 30 De Octubre De 2022];36(2):242-6. Disponible en: Https://Revistas.Um.Es/Analesps/Article/View/350301

[20] Buela G, Miró E, Iañez M. Catena A. Relación Entre La Duración Habitual Del Sueño Y El Estado De Ánimo Depresivo: Síntomas Somáticos Versus Cognitivos. Revista Internacional De Psicología Clínica Y De La Salud [Internet]. 2007[Citado 30 De Octubre De 2022];7(3):615-631. Disponible en: Https://Www.Redalyc.Org/Articulo.Oa?Id=33770303

[21] Concha C, González G, Piñuñuri R, Valenzuela C. Relación Entre Tiempos De Alimentación, Composición Nutricional Del Desayuno Y Estado Nutricional En Estudiantes Universitarios De Valparaíso, Chile. Rev. Chil. Nutr. [Internet]. 2019 [Citado 30 De Octubre 2022 ] ; 46( 4 ): 400-408. Disponible En:Https://Www.Scielo.Cl/Scielo.Php?Pid=S0717-75182019000400400&Script=Sci_Arttext

[22] Carrión C, Zavala I. El Estado Nutricional Asociado a Los Hábitos Alimentarios Y El Nivel De Actividad Física De Los Estudiantes De La Facultad De Ciencias De La Salud De La Universidad Católica Sedes Sapientiae en el Periodo 2016 – II [Tesis T].Lima: Universidad Católica Sedes Sapientiae Lima Tesis 2016 [Citado 30 De Octubre 2022 ] ; Disponible En: https://repositorio.ucss.edu.pe/bitstream/handle/20.500.14095/547/Carrion_Zavala_tesis_bachiller_2018.pdf?sequence=1&isAllowed=y

[23] León S, Aníbal R, Guerrero M, Luis R. Estilo De Vida Y Salud. Educere [Internet]. 2010. [Citado 30 De Octubre 2022] ; 14( 48 ): 13-19. Disponible en: Https://Www.Redalyc.Org/Articulo.Oa?Id=35616720002

[24] Carranza E, Caycho T, Salinas S, Ramírez M, Campos C, Chuquista I K Et Al. Efectividad De Intervención Basada En Modelo De Nola Pender En Promoción De Estilos De Vida Saludables De Universitarios Peruanos. Rev Cub. Enfermer. [Internet]. 2019 [Citado 29 Octubre 2022 ] ; 35( 4 ): 28-59. Disponible en: .Http://Scielo.Sld.Cu/Scielo.Php?Script=Sci_Arttext&Pid=S0864-03192019000400009#B8

[25] Palomares L. "Estilos De Vida Saludables Y Su Relación Con El Estado Nutricional En Profesionales De La Salud". [Tesis De Grado De Magíster]. Lima: Universidad Peruana De Ciencias Aplicadas. Programa De Maestría En Gestión Y Docencia En Alimentación Y Nutrición; 2014. Disponible en: http://hdl.handle.net/10757/566985

[26] Arrieta A. Relación Entre Aptitud Física, Estado Nutricional Y Nivel De Actividad Física En Bomberos Pertenecientes A Compañías De Lima Y

Callao, 2020. [Tesis De Grado].Lima: Universidad Nacional Mayor De San Marcos Universidad Del Perú. Decana De América Facultad De Medicina Escuela Profesional De Nutrición; 2020. Disponible en: Http://Cybertesis.Unmsm.Edu.Pe/Bitstream/Handle/20.500.12672/14001/Arrieta_Aa.Pdf?Sequence=1&Isallowed=Y

[27] Organización Mundial De La Salud. Malnutrición. [Internet]. Ginebra Oms;2021[Citado 30 De Octubre 2022 ]. Disponible en :Https://Www.Who.Int/Es/News-Room/Fact-Sheets/Detail/Malnutrition

[28] Morales J, Hernán M, Et Al. Exceso De Peso Y Riesgo Cardiometabólico En Docentes De Una Universidad De Lima: Oportunidad Para Construir Entornos Saludables. Educ Med. [Internet]. 2018 [Citado El 30 De Octubre De 2022];19(S3):256-261. Disponible En: Https://Doi.Org/10.1016/J.Edumed.2017.08.003

[29] Castro J, Cerna I. Hábitos alimentarios, estado nutricional y riesgo cardiovascular en bomberos de 20 a 59 años del Batallón XII, Costa Rica, 2020.Rev Hisp Cienc Salud. [interne ] 2020 [citado 17 de noviembre de 2022]; 6(4):166-17. Disponible en: https://uhsalud.com/index.php/revhispano/article/view/446/277

[30] Echeverría M. Estado Nutricional Y Hábitos Alimentarios Del Personal De Cuerpo De Bomberos De Cantón Otavalo 2021. Ecuador [Tesis Previa La Obtención Del Título]. Ecuador; Universidad Técnica Del Norte Facultad Ciencias De La Salud Carrera De Nutrición Y Salud Comunitaria. 2021 [Citado El 30 De Octubre De 2022].Disponible En : http://repositorio.utn.edu.ec/handle/123456789/11141

[31] Camargo F, Jardim T, Rocha L, Zandonade E, Nívea K. Prevalencia de obesidad en bomberos brasileños y la asociación de la obesidad central con factores de riesgo personales, ocupacional y cardiovasculares: un estudio transversal. BMJ Open. [Internet]. 2020 [Citado en 10 de abril del 2023]; 12(10). Disponible en: https://bmjopen.bmj.com/content/bmjopen/10/3/e032933.full.pdf

[32] Chuquipoma J. Relación Entre Los Conocimientos Previos En Nutrición Y El Estado Nutricional En Bomberos De La Compañía "Salvadora Lima N° 10", 2018" [Tesis Para Optar El Título Profesional De Licenciada En Nutrición]. Universidad Nacional Villarreal.2019 [Citado El 30 de octubre De 2022]. Disponible en: Https://Repositorio.Unfv.Edu.Pe/Bitstream/Handle/20.500.13084/2937/Unfv_Chuquipoma_%C3%91iquen_Jannet_Angelita_Titulo_Profesional_2019.Pdf?Sequence=1&Isallowed=Y

[33] Rodríguez C. Estilos De Vida Y Factores Biosocioculturales De Los Trabajadores Voluntarios De La Compañía De Bomberos B-107 Nuevo Chimbote, 2017" [Tesis Para Optar El Título Profesional De Licenciada En Enfermeria]. Universidad Católica Los Ángeles De Chimbote.2018 [Citado El 30 de Octubre de 2022]. Disponible en: Https://Repositorio.Uladech.Edu.Pe/Bitstream/Handle/20.500.13032/8564/Bomberos_Estilos_De_Vida_Rodriguez_Cabrera_Carolina_Marlit%20.Pdf?Sequence=1&Isallowed=Y

[34] Ministerio de Salud. Instituto Nacional de Salud. Guía técnica para la valoración nutricional antropométrica de la persona adulta [Internet]. Lima-Perú; 2012 [consultado 05 Noviembre 2022]. Disponible en: https://bvs.ins.gob.pe/insprint/CENAN/Valoraci%C3%B3n_nutricional_antropom%C3%A9trica_persona_adulta_mayor.pdf

[35] Gomez P. Principios básicos de bioética. Revista Peruana de Ginecología y Obstetricia [revista en Internet] 2009 [consultado 10 de Noviembre de 2022]; 55(4): 230-233. Disponible en: http://sisbib.unmsm.edu.pe/BVRevistas/ginecologia/vol55_n4/pdf/A03V55N4.pdf

[36] Carreño Dueñas J. Consentimiento informado en investigación clínica: Un proceso dinámico. Persona y Bioética [revista en Internet] 2016 [consultado 10 de noviembre de 2022]; 20(2): 232-243. Disponible en: http://personaybioetica.unisabana.edu.co/index.php/personaybioetica/article/view/232/html_1

[37] Arias S, Peñaranda F. La investigación éticamente reflexionada. Revista Facultad Nacional de Salud Pública [consultado 10 de noviembre de 2022]; 33(3): 444-451 [Internet]. 2015. Available from: http://www.scielo.org.co/scielo.php?script=sci_arttext&pid=S0120-386X2015000300015

# Enhanced Emotion Analysis Model using Machine Learning in Saudi Dialect: COVID-19 Vaccination Case Study

Abdulrahman O. Mostafa, Tarig M. Ahmed

Department of Information Technology-Faculty of Computing and Information Technology
King Abdulaziz University, Jeddah 21589, Saudi Arabia

*Abstract*—Sentiment Analysis (SA) and Emotion Analysis (EA) are effective areas of research aimed to auto-detect and recognize the sentiment expressed in a text and identify the underpinning opinion towards a specific topic. Although they are often considered interchangeable terms, they have slight differences. The primary purpose of SA is to find the polarity expressed in a text by distinguishing between positive, negative, and neutral opinions. EA is concerned with detecting more emotion categories, such as happiness, anger, sadness, and fear. EA allows the analysis to extract more accurate and detailed results that suit the field in which it is applied. This work delves into EA within the Saudi Arabian dialect, focusing on sentiments related to COVID-19 vaccination campaigns. Our endeavor addresses the absence of research on developing an effective EA machine-learning model for Saudi dialect texts, particularly within the healthcare and vaccinations domain, exacerbated by the lack of an EA manual-labeled corpus. Using a systematic approach, a dataset of 33,373 tweets is collected, annotated, and preprocessed. Thirty-six machine learning experiments encompassing SVM, Logistic Regression, Decision Tree models, three stemming techniques, and four feature extraction methods enhance the understanding of public sentiment surrounding COVID-19 vaccination campaigns. Our Logistic Regression model achieved 74.95% accuracy. Findings reveal a predominantly positive sentiment, particularly happiness, among Saudi citizens. This research contributes valuable insights for healthcare communication, public sentiment monitoring, and decision-making while providing labeled-corpus and ML model comparison results for improving model performance and exploring broader linguistic and dialectal applications.

*Keywords—Data mining; natural language processing; sentiment analysis; emotion analysis; machine learning; support vector machine; logistic regression; decision tree; Covid-19*

## I. INTRODUCTION

Microblogging has grown significantly as a means of communication and information sharing, notably on platforms like Twitter. Offering real-time accessibility from anywhere, Twitter allows users, through short texts or "Tweets," to express thoughts, opinions, and emotions [1]. This social media giant has become a global hub for diverse purposes, including news updates, social interactions, and discussions. In Saudi Arabia, where Twitter boasts over 15.5 million active users [2], it is a valuable repository for researchers keen on comprehending public sentiments.

The advent of Natural Language Processing (NLP) techniques, such as Sentiment Analysis (SA) and Emotion Analysis (EA), has revolutionized the understanding of extensive textual data generated on platforms like Twitter. SA discerns a text's sentiment or emotional tone, while EA goes a step further to identify specific emotions like happiness, anger, or sadness. These techniques find widespread applications in diverse domains, including social media analytics [5], customer feedback analysis [6], and market research, offering insights into human opinions, emotions, and behaviors.

The research gap addressed in this study emerges from several significant factors within the context of emotion analysis in the Saudi dialect. Arabic is known for its rich morphology and many dialects [3]. Social media, particularly Twitter, introduces informal language, including dialects and slang, complicating the accurate interpretation of emotions [2]. Additionally, the absence of a multi-emotion class Saudi dialect labeled-tweets corpus, the presence of diacritical marks (Tashkeel) introducing ambiguity, and the deviation of the Saudi dialect from Modern Standard Arabic writing norms collectively contribute to this research gap. The study's overarching objective is to fill this void by crafting a machine-learning model adept at classifying Saudi dialect tweets into seven emotion categories, explicitly focusing on emotions related to COVID-19 vaccinations. Such an endeavor is poised to enhance our understanding of the sentiments expressed in Saudi Arabia concerning COVID-19 vaccinations.

Existing studies have provided limited insights using narrow classifications, making a more comprehensive range of classifications necessary for more valuable and effective results. Embracing EA offers a more comprehensive understanding of sentiments beyond binary or ternary classifications, particularly in contexts like COVID-19 vaccination campaigns.

The motivations behind the study are twofold. First, there is a need to focus on implementing an EA model in the Saudi dialect, anticipating its profound influence across various sectors such as business, healthcare, education, government, and technology. Second, to better understand the general attitudes towards COVID-19 vaccination campaigns held in Saudi Arabia. Objectives are aligned with these motivations, striving to produce an effective machine-learning model, enhance the accuracy of existing EA studies, unveil prevailing

attitudes, create a Saudi dialect labeled-tweets corpus, and evaluate different algorithms comprehensively.

In anticipation of these objectives, the study expects to deliver a machine-learning model proficiently classifying Saudi dialect tweets into seven emotion categories. Additionally, it aims to create a labeled-tweets corpus, shedding light on the emotions expressed in diverse contexts, particularly in the healthcare and COVID-19 vaccination domain. Visual representations of general attitudes and detailed statistics are poised to empower decision-makers with nuanced insights, influencing policies and communication strategies. The model and corpus crafted in this study are envisioned to be valuable assets, not only for this research domain but also for broader applications in related studies.

Our paper is structured as follows: In Section II, we provide an overview of sentiment analysis and emotion analysis, highlighting their intersections and briefly discussing prior academic work in the field. Section III details the methodology, materials, and steps in constructing the final machine learning model, covering the dataset collection, annotation, and balancing methods. Section IV outlines the implementation of three machine learning models. Section V covers the evaluation methods and discusses the results achieved. Finally, in Section VI, we conclude the paper, presenting a summary and outlining our future work.

## II. BACKGROUND

### A. Sentiment Analysis

*1) Using machine learning:* In health-related contributions addressing the COVID-19 pandemic, Aljameel et al. [18] developed a machine learning model gauging individuals' awareness of preventive measures during the quarantine in Saudi Arabia. Utilizing a dataset from the curfew period, they employed SVM, NB, and KNN classifiers, optimizing 85% accuracy by combining TF-IDF with SVM. Focused on the Saudi dialect, the study by Al Sari et al. [19] in the entertainment field used MLP, NB, SVM, RF, and Voting algorithms, achieving 90% accuracy with NB and MLP on Twitter data. Alhuri et al. [20] utilized GRU in an RNN, reaching an 81% F1 score for public reactions to COVID-19 in Arabic tweets. Alahmary, Al-Dossari, and Emam. The study in [4] outperformed ML with DL algorithms (Bi-LSTM) on the Saudi Dialect Twitter Corpus, achieving 94% accuracy. AlYami and AlZaidy [21] focused on Arabic Dialect Identification using SVM, RF, NB, and LR, achieving a maximum of 87% and 86% accuracy in Egyptian dialects using LR and SVM, respectively. However, they have worked on sentiment analysis fields only; limitations included a small dataset, the absence of manual annotation, and intensive preprocessing.

*2) Using lexicon-based:* Assiri, Emam, and Al-Dossari [3] proposed a domain-specific lexicon-based algorithm for the Saudi dialect, addressing the absence of such models. However, limitations arose from evaluating it against a non-Saudi dataset and focusing solely on text polarity. Similarly, Al-Thubaity et al. [10] manually created "SauDiSenti," a

Saudi dialect sentiment lexicon, but challenges emerged in comparing it to a broader Arabic dictionary. Al-Ghaith [22] adopted a distinct approach, enhancing sentiment analysis accuracy by directly applying preprocessing tasks to the Saudi dialect lexicon, achieving 81% accuracy by relying on an English lexicon for the original creation.

*3) Using hybrid approach:* Very few published studies have utilized the Hybrid approach of SA in Arabic - where semantic orientation and ML techniques are combined. Aldayel and Azmi [9] used this approach to improve the F-measure score; they achieved an overall F-measure and accuracy of 84% and 84.01%, respectively. Alhumoud, Albuhairi, and Altuwaijri [23] used the same approach of combining two ML algorithms and applied the SA on 3000 Saudi dialect tweets to prove the efficiency of the hybrid learning approach compared to solo ML.

### B. Emotion Analysis

EA can be described as recognizing distinct human emotions in contrast to Sentiment Analysis, which identifies whether data is positive, negative, or neutral [24]. Because of the lack of a labeled corpus for Saudi dialect that can be used for classifying emotions and polarity behind a text, Al-Thubaity, Alharbi, Alqahtani and Aljandal [10] introduced the Saudi Dialect Twitter Corpus (SDTC) that contains 5400 tweets of Saudi dialect and MSA classified for sentiment analysis and emotion analysis annotated by three raters based on their polarity (positive, negative, and neutral) for the sentiment, and based on Ekman's basic emotions (anger, fear, disgust, sadness, happiness, surprise, no emotion and not sure) for the emotion analysis. However, no ML or lexicon-based approaches have been applied to this corpus to evaluate its efficiency in this study. Another study by A. AlFutamani and H. Al-Baity [11] was the first in the EA in Arabic, focusing on Saudi dialects in Arabic textual content retrieved from Twitter, mainly in Saudi-based tweets. They built a system that can detect the underlying emotions of Saudi dialect tweets to classify them based on seven emotion categories (happiness, fear, disgust, anger, surprise, optimism, and sadness). They used two ML algorithms (SVM and MNB), achieving 73.39% accuracy in the SVM approach. However, they applied the analysis in dataset domain sets different from ours with varying models of ML.

In summarizing the related work, it is evident that sentiment analysis and emotion analysis have made significant strides, particularly in applications related to COVID-19 discussions on social media. However, within the context of our study, there exists a notable research gap. Previous works have primarily addressed sentiment analysis in broader contexts, lacking the depth to decipher the intricate emotional expressions specific to the Saudi population. Furthermore, the scarcity of labeled datasets in the Saudi dialect poses a considerable challenge. Our research seeks to bridge this gap by offering a comprehensive analysis of emotion analysis, utilizing a meticulously annotated dataset tailored to the Saudi dialect. This approach contributes to the broader field of sentiment analysis and provides a nuanced understanding of

the emotional landscape surrounding COVID-19 vaccinations in Saudi Arabia.

### C. Machine Learning Algorithms

ML algorithms encompass supervised, unsupervised, and reinforcement learning. Supervised learning trains on labeled data, associating input with output labels. It excels in classification, categorizing data into known classes, regression for predicting continuous values, and ranking for ordering data [7] [8]. Techniques like Support Vector Machines, Logistic Regression, Decision Trees, Random Forests, and Neural Networks are part of supervised learning, each chosen based on factors like data nature and problem complexity. Unsupervised learning handles unlabeled data, discerning patterns without explicit guidance, while reinforcement learning involves agent learning decisions to maximize rewards in an environment [12] [13]. Supervised learning's versatility finds applications in various domains, offering solutions in classification, regression, and ranking [14] [15].

*1) Support Vector Machine (SVM):* Stands out as a prominent supervised learning algorithm, initially crafted by Vapnik for binary classification and regression [37]. Renowned for its robust theoretical foundations, SVM has evolved to address multi-class classification using techniques like one-vs-rest and one-vs-one approaches. In the one-vs-rest method, distinct SVM models are trained for each class, treating it as positive and the others as negative, with the final result determined by the most probable classifier. Conversely, the one-vs-one strategy involves SVM models comparing each class against every other class, employing a voting scheme for the outcome. SVM excels in handling both linearly and non-linearly separable data. For linearly separable data, the hyperplane equation is defined as $g(x) = w^T x + b$, where w is the weight vector, x is the input data vector, and b is the bias term. The Euclidean norm determines the vector's magnitude, which is crucial for understanding its length in n-dimensional space. The optimization goal of SVM is to find the optimal hyperplane, maximizing the margin between support vectors accomplished through convex quadratic programming. SVM employs the kernel trick for non-linearly separable data, transforming input data into a higher-dimensional feature space for linear separation. The choice of the kernel function, whether linear, polynomial, radial basis function (RBF) or HyperTangent, profoundly influences SVM's ability to capture intricate patterns and relationships in the data [38]. SVM's suitability for text classification stems from its generalization capabilities, adherence to the Structural Risk Minimization principle, capacity to incorporate prior knowledge, and superior performance to alternatives like k-nearest-neighbors (kNN).

*2) Logistic* regression (LR) is a widely employed supervised learning algorithm that is a statistical method for modeling the probability of a binary outcome based on predictor variables. LR recognizes vectors with variables in text classification, assesses coefficients for each input, and predicts text class as a word vector. This model measures the statistical significance of independent variables concerning probability, offering a potent means of modeling binomial outcomes. LR excels in text categorization, providing advantages like computing probability values instead of scores. The logistic function, or sigmoid function, characterizes the LR model's relationship between variables and the probability of the outcome. Ensuring predicted probabilities fall within the range of 0 and 1, the LR equation is $P = 1/(1+e^{(-(w+bX))})$. In training, LR estimates parameters (weights) w and b through maximum likelihood estimation, aiming to maximize the likelihood of observed data. LR is computationally efficient, easy to implement, and yields interpretable results with estimated coefficients. It accommodates numerical and categorical input features and extends to multi-class classification using strategies like one-vs-rest, where a separate LR model is trained for each class, determining the final prediction based on the highest probability among all models.

*3) Decision Tree (DT):* A machine-learning algorithm for classification and regression tasks. It operates by recursively partitioning the data based on the values of input features, ultimately leading to a decision regarding the target variable. The algorithm constructs a tree-like structure representing a sequence of decisions and their potential consequences. Each internal node of the tree corresponds to a test on a specific attribute, while each branch represents the outcome of the test. The tree's leaf nodes correspond to class labels or numerical values [16]. Decision trees are popular in machine learning due to their interpretability, simplicity, and ability to handle various data types, including categorical and numerical variables [17]. They find applications in multiple domains, including image processing, clinical practice, and financial analysis. Their inherent structure allows for an intuitive understanding of the decision-making process, making them valuable tools for extracting insights from data.

### D. Covid-19 and Vaccinations

The COVID-19 pandemic, caused by the novel coronavirus SARS-CoV-2, has profoundly impacted societies worldwide [28]. It was first identified when the initial case was reported in Wuhan, China, on December 19, 2019 [29]. The World Health Organization (WHO) officially declared the global COVID-19 pandemic on March 11, 2020 [30]. This declaration marked a turning point in the international response to the virus, leading to widespread public health measures to curb its spread. Saudi Arabia recorded its first confirmed case of COVID-19 on March 2, 2020 [31] [32]. The country swiftly implemented various measures to combat the virus's transmission, including lockdowns and travel restrictions.

The introduction of COVID-19 vaccines marked a pivotal moment in the fight against the pandemic. Saudi Arabia approved the Pfizer-BioNTech vaccine on December 10, 2020 [33]. Registration for the vaccine in Saudi Arabia began on December 15, 2020 [34]. As vaccination efforts progressed, restrictions evolved, with announcements such as allowing only vaccinated individuals to enter certain buildings starting from August 1, 2021. The vaccine rollout continued to

advance, with the commencement of the second vaccine dose administration on June 23, 2021 [35]. As the situation improved, Saudi Arabia took steps to return to normalcy, including lifting many precautionary measures on May 3, 2022 [36].

### III. MATERIAL AND METHOD

This section provides a comprehensive overview of our methodology for producing a machine learning (ML) model and constructing a labeled corpus for emotion analysis in healthcare, specifically focusing on the Saudi dialect. The process comprises distinct stages, outlined in Fig. 1. The initial step involves collecting relevant tweets on vaccinations through specific keywords using the Twitter API. These tweets are then manually annotated by three Saudi natives, forming a labeled corpus for training and evaluation. Data preprocessing is undertaken to filter irrelevant content, rectify text errors, ensure consistency, and tokenize the text for analysis. Feature extraction follows, where we employ Bag-of-Words, N-Gram, and TF-IDF methods to effectively capture emotions expressed in the tweets, serving as inputs for the ML model. The classification stage involves developing and training the ML model, utilizing Support Vector Machines, Logistic Regression, and Decision trees to classify tweets into seven emotions. The performance evaluation includes using various metrics and experiments to validate the methodology. Unseen data is employed for testing, with metrics like accuracy, precision, recall, and F1-score measured to assess the effectiveness and limitations of our approach.



Fig. 1. Block diagram of the research methodology.

#### A. Dataset Collection

The data collection phase is pivotal to our implementation, aiming to amass diverse tweets related to vaccinations in the Saudi dialect. We utilized the Twitter API for Academic Research with Python programming language to access real-time Twitter data, employing libraries such as Tweepy [39], Pandas, and Requests. The collection spanned from December 15, 2020, to March 10, 2022, crucial periods in Saudi Arabia's COVID-19 vaccination timeline. Specific Arabic keywords targeting vaccination-related discussions were specified in each request, ensuring relevance. Additional parameters refined the

dataset, focusing on the Saudi region and Arabic language and excluding retweets. The process yielded 34,074 raw tweets, each characterized by properties like tweet text, author ID, creation time, geolocation, and engagement metrics. The dataset, saved in CSV format, forms a robust foundation for subsequent annotation, preprocessing, feature extraction, and classification stages. Table I shows one raw instance of collected tweets.

TABLE I. ONE RAW INSTANCE OF THE COLLECTED TWEETS

| Tweet | انتهت حفلة كورونا!! وقطاع السوبرماركت @FBasmer @_doje_ |
|---|---|
| Author_ID | 534798407 |
| Created_at | 2020-12-19 17:37:51+00:00 |
| Geo | 000799c66e428a87 |
| ID | 1.34035E+18 |
| Lang | Ar |
| Like_count | 0 |
| Quote_count | 0 |
| Reply_count | 0 |
| Retweet_count | 0 |
| Source | Twitter for Android |

#### B. Data Annotation

In the data annotation phase, emotions were manually assigned to collected tweets following initial preprocessing steps. These steps included eliminating duplicates, Twitter handles, URLs, English text, and emojis/non-emoji symbols. The annotation categorized tweets into emotion classes aligned with Paul Ekman's emotions [40], augmented by a neutral/spam class and an optimistic class to account for Saudi dialect nuances. Eight emotion categories were established, detailed in Table II.

TABLE II. SHOWS THE EMOTIONS USED IN THE ANNOTATION STAGE

| # | English Emotion Class | Arabic Emotion Class | Explanation |
|---|---|---|---|
| 1 | happiness | سعادة | When the tweet expresses feelings of joy, happiness, or delight |
| 2 | fear | خوف | When the tweet expresses feelings of fear |
| 3 | disgust | اشمئزاز | when the tweet expresses disgust, disgust, or disgust with the tweeter |
| 4 | anger | غضب | When the tweet expresses angry feelings |
| 5 | surprise | تفاجؤ | When the tweet expresses feelings of surprise, astonishment, or wonder |
| 6 | optimism | تفاؤل | When the tweet expresses feelings of optimism and a positive view of the future |
| 7 | sadness | حزن | When the tweet expresses feelings of sadness, brokenness, grief, or depression |
| 8 | neutral | محايد/إعلان/لا يمكن تحديده | When the content of the tweet does not include any emotions or feelings that cannot be identified from other options or includes only advertising hashtags without any emotions |

Google Sheets were employed in the annotation stage, where a copy of tweets was shared with individual annotators familiar with the Saudi dialect. Annotators were guided by detailed instructions ensuring privacy, single emotion selection for ambiguous cases, accuracy, and time management. The annotation process lasted two months, maintaining a manageable daily average for annotators. Results were consolidated into a dataset, revealing 23,689 fully matched tweets and 9,684 with discrepancies. Table III illustrates counts for each emotion class by annotators, and Fig. 2 provides examples of annotated tweets. This annotation phase produced a labeled dataset, a crucial foundation for subsequent stages like data preprocessing, feature extraction, and classification.

*C. Preprocessing*

The preprocessing stage, integral to our methodology, comprised two pivotal sub-stages: pre-annotation and post-annotation. In the pre-annotation phase, executed within Google Sheets using sophisticated REGEX formulas, we undertook various measures to refine the dataset comprehensively. Initially, we focused on eliminating redundancy, ensuring uniqueness in our dataset by removing 701 duplicate tweets from the initial count of 34,074, resulting in 33,373 distinct tweets. Furthermore, we removed Twitter handles (@), URLs, English text, emojis, and non-emoji symbols to enhance the dataset's purity. The cleaning tasks aimed at centering our analysis on the core content of tweets, devoid of any external influences. In the post-annotation, we employed the KNIME [41] platform to undertake further profound preprocessing exclusively on the subset of class-matched tweets from the three raters, totaling 22,689 tweets. Fig. 3 demonstrates the comprehensive nature of our data refinement process.

The post-annotation preprocessing involved a multifaceted approach:

- Punctuation Removal: Punctuation marks were expunged from tweet text to eliminate unnecessary noise, allowing focused analysis of words and their emotional significance.

- Elimination of Numbers: Numeric characters were systematically removed from tweets as they do not contribute directly to emotional content. This step simplified the text and heightened subsequent analysis accuracy.

TABLE III. COUNTS TWEETS ANNOTATED IN EACH EMOTION CLASS BY THREE ANNOTATORS

| Rater | | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|---|---|---|---|
| 1# | # | 21992 | 1723 | 1169 | 756 | 1358 | 821 | 1277 | 4277 |
| | % | 65.90% | 5.16% | 3.50% | 2.27% | 4.06% | 2.47% | 3.82% | 12.82% |
| 2# | # | 26491 | 794 | 1177 | 19 | 1723 | 179 | 576 | 2414 |
| | % | 79.38% | 2.38% | 3.53% | 0.06% | 5.16% | 0.53% | 1.73% | 7.23% |
| 3# | # | 22415 | 1313 | 1069 | 419 | 1499 | 840 | 1214 | 4604 |
| | % | 67.16% | 3.93% | 3.20% | 1.25% | 4.49% | 2.51% | 3.63% | 13.79% |

| Tweet | class by rater #1 | class by rater #2 | class by rater #3 |
|---|---|---|---|
| مدينة الضباب تغلق ابوابها ابتداء من الغد ديسمبر بسبب السلاله الجديده المتجوره من بروس كورونا اغلاق تام للبوتيكات والمطاعم والصالونات وجميع المحلات الغير اساسيه فقط الصيدليات والسوبر ماركات بتكون فاتحه | 5 (surprise) | 8 (neutral/spam) | 1 (happiness) |
| الحمداله اخذت اول جرعة من اللقاح اللهم انفع بها خذ الخطوة | 1 (happiness) | 1 (happiness) | 1 (happiness) |
| تعبت نفسيا والسبب ان كورونا للحين مو راضي يخلص | 7 (sadness) | 8 (neutral/spam) | 2 (fearness) |
| اذا شركة واحدة اللي يتصنع اللقاح بيكون احتكار والشركات الطبية كل سنة تزيد توسعها علشان الدواء مثل له اكثر من شركة مصنعة هو نفس الحال مع كورونا | 8 (neutral/spam) | 8 (neutral/spam) | 8 (neutral/spam) |

Fig. 2. Examples of annotated tweets by different raters.



Fig. 3. The flow of preprocessing nodes in the KNIME platform.

- Removal of Double Spaces and New Lines: Consecutive spaces and new lines were eradicated, ensuring uniformity in text format.

- Arabic Diacritics Removal and Character Normalization: Diacritics, such as vowel marks, were deleted for consistency. Characters were normalized to ensure uniformity and standardization across the dataset. This eliminated text variations that could impact emotion analysis accuracy, as shown in Fig. 4.



Fig. 4. An example of a tweet before and after diacritics removal.

- Normalization Techniques: Normalization techniques were applied to standardize data and convert characters to their original shape. For instance, sequences of certain Arabic letters were regulated for uniformity.

- Stemming: Employing three distinct stemming techniques (Snowball Stemmer, Porter Stemmer, and Kuhlen Stemmer), words were reduced to their root form. This facilitated a more accurate analysis and interpretation of emotions.

- Stopwords that lack significant meaning were removed using a Stop Word Filter node in KNIME. A stop-word dictionary for the Arabic language was employed for this purpose [42].

- Tokenization: Text was tokenized into individual words using the Arabic tokenizer provided by the NLTK library.

- Column Filtering: Irrelevant columns were filtered out, retaining only tweet and label class columns, streamlining subsequent emotion analysis.

This meticulous preprocessing led to a pristine dataset comprising 22,549 clean tweets. This refined dataset is now poised for the subsequent stages of feature extraction, classification, and further analysis, as shown in Fig. 5.



Fig. 5. Counts the tweets at each stage of the implementation.

### D. Feature Extraction and Selection

In the feature extraction and selection stage, our initial step involved the removal of the 8th emotion class label, encompassing neutral, spam, and advertisement text (totaling 19,197 tweets). This exclusion aimed to streamline subsequent analysis, focusing solely on the emotions inherent in the remaining dataset of 3,352 tweets.

*1) Bag-of-Words (BoW):* The Bag-of-Words approach, a fundamental yet robust technique, was employed to convert preprocessed text data into numerical vectors. Each word in the text is treated as a distinct feature, and its frequency is quantified in this method.

The resulting vector encapsulates the occurrence of each word in the text, irrespective of word order or grammar. Two sub-flows of BoW were applied to evaluate model accuracy, each depicted in. This method yielded a high-dimensional representation of the text data, forming the input for subsequent machine-learning models. Fig. 6 and Fig. 7 show these flows.



Fig. 6. First sequence of BoW flow.

Fig. 7.   Second sequence of BoW flow.

*2) N-Gram:* Extending the Bag-of-Words approach, the N-Gram method considers sequences of N consecutive words as features. We specifically employed bi-grams (N=2), representing pairs of successive words in each tweet. This extension allows for capturing context and word relationships in the text, enhancing our ability to discern nuanced insights into expressed emotions.

*3) Term Frequency-Inverse Document Frequency (TF-IDF):* TF-IDF, a prevalent technique in natural language processing, gauges the importance of each word in a document relative to a corpus. It factors in the frequency of a word in a specific document and its rarity across the entire corpus. Assigning higher weights to words that are frequent in a document but rare across the corpus makes them more discriminative. TF-IDF was employed to enrich the representation of words based on their significance in each tweet and across the entire dataset.

We transformed the preprocessed text data into numerical representations using these three feature extraction methods. These representations effectively captured essential information about the emotions expressed in the tweets. The resulting feature matrices served as inputs for the machine learning model during training, facilitating the model's ability to discern patterns and relationships between features and labeled emotions.

*E.  Resolve Data Imbalance (Oversampling)*

The focus is on addressing class imbalance, a crucial consideration in developing an effective emotion analysis model. The section outlines the strategies employed, particularly data partitioning and oversampling techniques.

*1) Class distribution:* An evaluation of the initial class distribution within the dataset is conducted before diving into oversampling. It's revealed that the original dataset of 3,352 tweets exhibits an imbalance across emotion categories, a factor that can impact the model's ability to discern less frequent emotions effectively as shown in Table IV and Fig. 8.

TABLE IV.     CLASS DISTRIBUTIONS OF THE 3352 TWEETS BEFORE DATA SPLITTING

| Class | Tweets Count |
|---|---|
| Happiness (1) | 1,926 |
| Fear (2) | 230 |
| Disgust (3) | 71 |
| Anger (4) | 482 |
| Surprise (5) | 13 |
| Optimism (6) | 243 |
| Sadness (7) | 387 |
| **Total** | **3352** |



Fig. 8.   Class distributions of the 3352 tweets before data splitting.

Fig. 9. Class distribution of the training dataset contains 2346 tweets after data splitting and before oversampling.

*2) Dataset splitting:* A 70:30 train-test split ratio is employed to evaluate machine learning model performance. This ensures that 70% of the data is allocated for training, and the remaining 30% is reserved for evaluation. The class distribution post-splitting is presented for both the training and testing datasets as shown in Table V and Fig. 9.

TABLE V. THE CLASS DISTRIBUTION OF THE TRAINING DATASET CONTAINS 2346 TWEETS AFTER DATA SPLITTING AND BEFORE OVERSAMPLING

| Class | Tweets Count |
|---|---|
| Happiness (1) | 1,348 |
| Fear (2) | 161 |
| Disgust (3) | 50 |
| Anger (4) | 337 |
| Surprise (5) | 9 |
| Optimism (6) | 170 |
| Sadness (7) | 271 |
| **Total** | **2346** |

*3) Oversampling:* Oversampling techniques are strategically applied to rectify the class imbalance within the training data. The oversampling is exclusively directed at the training data. After this process, the class distribution in the training dataset is significantly altered, as illustrated in Table VI and Fig. 10.

TABLE VI. CLASS DISTRIBUTION OF THE TRAINING DATASET CONTAINS 9399 TWEETS AFTER DATA SPLITTING AND AFTER OVERSAMPLING

| Class | Tweets Count |
|---|---|
| Happiness (1) | 1,348 |
| Fear (2) | 1288 |
| Disgust (3) | 1350 |
| Anger (4) | 1348 |
| Surprise (5) | 1350 |
| Optimism (6) | 1360 |
| Sadness (7) | 1355 |
| **Total** | **9399** |



Fig. 10. Class distribution of the training dataset contains 9399 tweets after data splitting and after oversampling.

This oversampling endeavor aims to equalize the representation of emotions across classes in the training dataset, providing a more balanced learning experience for machine learning models.

## IV. IMPLEMENTATION

In the Classification stage, we embarked on the crucial task of predicting the emotions expressed in the 9399 tweets using the extracted features. Indeed, the combination of three stemming techniques, four feature extraction methods, and three machine learning algorithms resulted in 36 different models being trained in this stage. Each model represents a unique configuration of the preprocessing and classification pipeline, contributing to the comprehensive evaluation and comparison of various approaches. We obtained an enriched dataset with numerical representations of the extracted features. However, the emotion labels must be converted to numerical classes to train the machine learning models. This transformation involved mapping each emotion category to a unique numerical class, making the dataset suitable for classification. Having prepared the data and set up the train-test split, we applied three widely used machine learning algorithms: Support Vector Machine (SVM), Logistic Regression, and Decision Tree. Each algorithm underwent training on the training data to learn the underlying patterns and relationships between the extracted features and emotions. After training, the models were tested using the test data to evaluate their predictive capabilities.

In the learner node of each machine learning algorithm, specific configurations and options were carefully selected to optimize the performance of the models. For the Support Vector Machine (SVM), we utilized the polynomial kernel with power approximately equal to 1, bias around 1, and gamma set to around 1. The overlapping penalty was established within the range of 0.1 to 1, enabling us to control the influence of overlapping data points on the model's decision boundaries. In the case of logistic regression, we opted for the stochastic average gradient (SAG) solver, known for its efficiency [25] and ability to handle large datasets. The maximal number of epochs was set between 10 and 20, ensuring the algorithm converged to the optimal solution while avoiding overfitting. For the Decision Tree algorithm, we set the maximum number of patterns the tree will store to support highlighting to a default value of 10,000. This configuration allowed us to manage the complexity of the tree while still capturing the relevant patterns and relationships in the data flows illustrated.

## V. EVALUATION AND RESULTS

In this section, we will evaluate the model results using different metrics. We utilized the KNIME platform for the evaluation process, employing the Model Predictor and Scorer nodes. The Model Predictor node takes the test data partition from the partitioning node and the trained model from the learner node as inputs. It then uses the trained model to predict the emotion labels for the given test data. The Scorer node compares the predicted and actual labels, generating a confusion matrix displaying correct and incorrect predictions for each emotion class.

### A. Evaluation Methods

*1) Confusion matrix:* a fundamental evaluation tool that provides a tabular representation of the performance of a machine-learning classification model [26]. It gives a comparison between actual and predicted values as it displays the number of true positive (TP), true negative (TN), false positive (FP), and false negative (FN) predictions for each emotion class, which will be used in the measure. Furthermore, it allows us to assess the model's accuracy and ability to classify emotions correctly. The confusion matrix is a square matrix of size N x N, where N represents the number of emotion classes. We have a 7 x 7 confusion matrix in our specific case, as illustrated in Table VII.

TABLE VII. CONFUSION MATRIX FOR SEVEN CLASSES

| Actual Class | Happiness | Fear | Disgust | Anger | Sadness | Suprise | Optimism |
|---|---|---|---|---|---|---|---|
| Happiness | TP1 | FP1 | FP2 | FP3 | FP4 | FP5 | FP6 |
| Fear | FP7 | TP2 | FP8 | FP9 | FP10 | FP11 | FP12 |
| Disgust | FP13 | FP14 | TP3 | FP15 | FP16 | FP17 | FP18 |
| Anger | FP19 | FP20 | FP21 | TP4 | FP22 | FP23 | FP24 |
| Sadness | FP25 | FP26 | FP27 | FP28 | TP5 | FP29 | FP30 |
| Suprise | FP31 | FP32 | FP33 | FP34 | FP35 | TP6 | FP36 |
| Optimism | FP37 | FP38 | FP39 | FP40 | FP41 | FP42 | TP7 |

*2) Class distribution consideration:* As our dataset exhibits an imbalanced class distribution, we applied stratified cross-validation with some emotion classes having significantly fewer instances than others. Stratified cross-validation ensures that each fold retains the same proportion of instances for each emotion class as the original dataset. This approach is crucial for preventing biased evaluations and ensuring that each emotion class is represented appropriately during model training and testing.

*3) Comparative analysis:* To conduct a comparative analysis, we evaluated a total of 36 different models, considering the combination of three stemming techniques (Kohlen Stemmer, Porter Stemmer, and Snowball Stemmer), four feature extraction methods (Bag-of-Words, N-Gram, and TF-IDF), and three machine learning algorithms (Support Vector Machine, Logistic Regression, and Decision Tree). By comparing the performance of these models, we can identify the most practical combination of techniques and algorithms for sentiment analysis in the context of our research.

### B. Evaluation of Performance Metrics

A comprehensive set of performance metrics is employed to evaluate the effectiveness of emotion classification models for Saudi dialect tweets related to vaccinations. The chosen metrics, including Accuracy, Precision, Recall, and F1-score, offer valuable insights into the models' ability to classify emotions accurately.

*1) Accuracy:* a widely used metric in classification models, measures the proportion of correctly classified instances over the total number of instances in the dataset. It provides an overall assessment of the model's performance in correctly identifying emotions in tweets across all emotion classes. The formula for accuracy involves True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN).

*2) Precision:* assesses the proportion of true positive predictions for a specific emotion class over the total number of instances predicted as that class. It signifies the model's reliability in avoiding false positive predictions for a given class. Precision is crucial in understanding the model's accuracy when classifying tweets as a particular emotion.

*3) Recall:* sensitivity, or true positive rate, measures the proportion of true positive predictions for a particular emotion class over the total number of instances belonging to that class. It indicates the model's effectiveness in correctly identifying all instances of a specific emotion class. Recall is significant in understanding how well the model captures tweet emotions.

*4) F1-score:* a harmonic mean of precision and recall that combines both metrics into a single value. This score is instrumental in scenarios with an imbalance in class distribution. It offers a balanced evaluation of the model's performance, considering trade-offs between precision and recall in emotion classification.

*C. Results and Evaluation*

This section presents the outcomes of the 36 experiments conducted during the implementation stage. Table VIII, IX and Table X provide a comprehensive breakdown of the results for all 36 models, giving the performance metrics for each emotion class. These tables offer a basis for comparing each model's accuracy, precision, recall, and F1-score of each model across different feature extraction methods and stemming techniques. Among the configurations, the Logistic Regression model achieved the highest accuracy, reaching 74.95% when combined with N-Gram feature extraction and Snowball stemming; it also performs the same 74.95% accuracy when combined with the BoW with Snowball stemmer, followed by 74.75% Logistic Regression when combined with TF-IDF and Snowball Stemmer. The SVM model showcased a close performance with 74.35% accuracy when combined with N-Gram and Snowball.

Moving on to the recall metric, the SVM model achieved a remarkable 91.34% in both experiments, compromising N-Gram and TF-IDF when combined with the Snowball stemmer, displaying a higher ability to identify true positive instances correctly and making them top performers in emotion classification. The Logistic Regression result shows a close percentage of 91.00% with TF-IDF and Snowball stemmer.

Regarding precision, the SVM model demonstrated an impressive 98.98% precision rate when trained after using the N-Gram technique and Porter stemmer, indicating its capability to limit the number of false positives and ensure the accuracy of positive predictions. The Logistic Regression result shows a

close percentage of 97.23% with BoW and Kuhlen Stemmer, followed by another Logistic Regression experiment with Bow and Porter Stemmer achieving 96.60%.

The F-measure results highlighted the Logistic Regression model as the most balanced performer between precision and recall, achieving 92.93% and 92.51%, particularly when combined with TF-IDF and BoW feature extraction techniques with the Snowball stemmer. The SVM model followed closely in third place, achieving 92.38% when N-Gram and Snowball were used.

TABLE VIII. SVM MODEL RESULTS ARE BASED ON FOUR FEATURE EXTRACTION TECHNIQUES AND THREE STEMMING METHODS

| | Stemming | Accuracy % | Recall % | Precision % | F-Measure % |
|---|---|---|---|---|---|
| BoW-1 | Kuhlen | 64.31 | 83.04 | 95.61 | 88.88 |
| | Snowball | 67.39 | 85.46 | 96.10 | 90.47 |
| | Porter | 64.61 | 84.94 | 95.71 | 90.00 |
| BoW-2 | Kuhlen | 70.67 | 87.54 | 93.35 | 90.35 |
| | Snowball | 74.15 | 89.96 | 94.37 | 92.11 |
| | Porter | 70.57 | 87.02 | 93.32 | 90.06 |
| N-Gram | Kuhlen | 70.67 | 87.19 | 92.98 | 90.00 |
| | Snowball | 74.35 | 91.34 | 93.45 | 92.38 |
| | Porter | 70.17 | 87.19 | 98.98 | 90.00 |
| TF-IDF | Kuhlen | 70.17 | 87.19 | 92.98 | 90.00 |
| | Snowball | 74.35 | 91.34 | 93.45 | 92.38 |
| | Porter | 70.17 | 87.19 | 92.98 | 90.00 |

TABLE IX. DECISION TREE MODEL RESULTS ARE BASED ON FOUR FEATURE EXTRACTION TECHNIQUES AND THREE STEMMING METHODS

| | Stemming | Accuracy % | Recall % | Precision % | F-Measure % |
|---|---|---|---|---|---|
| BoW-1 | Kuhlen | 67.29 | 83.73 | 95.46 | 89.21 |
| | Snowball | 66.10 | 84.25 | 92.40 | 88.14 |
| | Porter | 67.29 | 83.73 | 95.46 | 89.21 |
| BoW-2 | Kuhlen | 66.00 | 84.42 | 94.39 | 89.13 |
| | Snowball | 69.38 | 88.23 | 90.58 | 89.39 |
| | Porter | 66.79 | 86.85 | 91.27 | 89.00 |
| N-Gram | Kuhlen | 66.79 | 86.85 | 91.27 | 89.00 |
| | Snowball | 69.38 | 88.23 | 90.58 | 89.39 |
| | Porter | 66.79 | 86.85 | 91.27 | 89.00 |
| TF-IDF | Kuhlen | 65.90 | 84.25 | 95.49 | 89.52 |
| | Snowball | 67.99 | 85.81 | 92.36 | 88.96 |
| | Porter | 66.00 | 84.42 | 94.39 | 89.13 |

*D. Discussion*

The results from our comprehensive experiment illuminate the effectiveness of our proposed Emotion Analysis machine learning model in classifying emotions and feelings expressed in Tweets written in the Saudi dialect. Our study's primary aim

was to develop a precise and dependable model that enhances existing study results and aligns with the Saudi dialect's unique linguistic and cultural intricacies. As demonstrated in Fig. 11 to Fig. 14. The most accurate, precise recall and F1-score models are highlighted. The Logistic Regression model, as depicted, surpasses others in terms of Accuracy and F-Measure metrics, signifying a significant achievement. The SVM model displays the highest Recall (sensitivity) and Precision performance. The third model (Decision Tree) is considered out of the competition of the top three. This success contributes to the validation of our objectives.

TABLE X.    SHOWS THE LOGISTIC REGRESSION MODEL RESULTS BASED ON FOUR FEATURE EXTRACTION TECHNIQUES AND THREE STEMMING METHODS

| | **Stemming** | **Accuracy %** | **Recall %** | **Precision %** | **F-Measure %** |
|---|---|---|---|---|---|
| BoW-1 | **Kuhlen** | 63.81 | 79.23 | 97.23 | 87.32 |
| | **Snowball** | 69.28 | 87.19 | 94.73 | 69.28 |
| | **Porter** | 66.00 | 83.73 | 96.60 | 89.71 |
| BoW-2 | **Kuhlen** | 72.86 | 89.79 | 93.68 | 91.69 |
| | **Snowball** | 74.95 | 90.83 | 94.25 | 92.51 |
| | **Porter** | 73.36 | 89.61 | 93.84 | 91.68 |
| N-Gram | **Kuhlen** | 73.06 | 89.96 | 94.71 | 92.28 |
| | **Snowball** | 74.95 | 90.31 | 94.05 | 74.95 |
| | **Porter** | 73.26 | 89.10 | 94.49 | 91.71 |
| TF-IDF | **Kuhlen** | 72.16 | 87.88 | 94.24 | 90.95 |
| | **Snowball** | 74.75 | 91.00 | 94.94 | 92.93 |
| | **Porter** | 72.76 | 88.58 | 93.60 | 91.02 |



Fig. 11.  The top accuracy results achieved by the three models.

When compared to existing works in the domain, our progress is remarkable. For instance, considering a study that achieved 73.39% accuracy through an SVM approach, which is 1.56% lower than our Logistic Regression result, and they solely employed two machine learning algorithms (SVM and MNB) [11], our accomplishments are striking. Not only did we surpass this accuracy benchmark, but we also encompassed a

diverse array of machine-learning algorithms, leading to a more robust and comprehensive evaluation. This also aligns with our objectives.

The performance of our model was significantly influenced by the careful selection of stemming techniques and feature extraction methods. Results indicate that specific combinations, such as employing N-Gram feature extraction with Snowball stemming, yield the highest accuracy. This underscores the importance of selecting the correct machine learning algorithm and optimizing preprocessing and feature engineering stages to exploit the data's potential fully.



Fig. 12.  The top recall results achieved by the three models.



Fig. 13.  The top precision results achieved by the three models.



Fig. 14.  Shows the top F-Measure Results achieved by the three models.

### E. General Attitudes Towards COVID-19 Vaccinations in Saudi Arabia

The outcomes of our model shed light on prevalent concerns and sentiments regarding COVID-19 vaccination in Saudi Arabia. The precision of our emotion analysis allows us to extract insightful understandings from public sentiment, aiding decision-making in healthcare initiatives. The distribution of annotated tweets across emotion classes is detailed in Table XI. They are revealing a diverse emotional landscape. 'Happiness' dominates, followed by 'Anger,' 'Sadness,' and 'Optimism.' 'Disgust' is rare, and 'Surprise' is infrequent. Insights and Latent Dirichlet Allocation (LDA) [27] indicate a prevailing positive disposition toward the COVID-19 vaccination campaign in Saudi Arabia. As shown in Fig. 15, 57.46% of tweets expressed happiness followed by anger 14.38%, sadness 11.55%, and optimism 7.25%. Fearful emotions account for 6.86%, while disgust and surprise are 2.12% and 0.39%, respectively. These findings highlight contentment with vaccine availability, achieving our objective of providing valuable insights into general attitudes toward vaccinations.



Fig. 15. The distribution of annotated tweets across various emotion classes.

TABLE XI. NUMBER OF INSTANCES IN EACH EMOTION CLASS IN THE FINAL STAGE OF IMPLEMENTATION

| Class | Number of Instances | % |
|---|---|---|
| 1- happiness | 1926 | 57.5 |
| 2- fear | 230 | 6.9 |
| 3- disgust | 71 | 2.1 |
| 4- anger | 482 | 14.4 |
| 5- surprise | 13 | 0.4 |
| 6- optimism | 243 | 7.2 |
| 7- sadness | 387 | 11.5 |

### F. Saudi Dialect Labeled-Tweets Corpus Availability

In alignment with our objective to produce a Saudi dialect labeled-tweets corpus in the healthcare and COVID-19 vaccination domain, we are pleased to announce the availability of the "saudiEAR" repository on GitHub [43]. This repository contains a comprehensive collection of tweets, including both the original dataset collected and the preprocessed version. By making this corpus publicly accessible, we aim to contribute to the research community and facilitate advancements in sentiment analysis, emotion classification, and related fields. This open dataset enables researchers and practitioners to explore the intricacies of sentiment expression in the Saudi dialect, particularly in healthcare and COVID-19 vaccination discussions.

## VI. CONCLUSION AND FUTURE WORK

### A. Conclusion

In conclusion, this research addresses a notable gap in emotional analysis, focusing on the Saudi context amid COVID-19 vaccination discussions. The objective was to build an effective machine-learning model for classifying Saudi tweets into distinct emotions, a need prompted by the scarcity of such studies in the Saudi dialect and the absence of a suitably labeled corpus. The methodology involved a meticulous collection of 34,074 Arabic tweets emphasizing COVID-19 vaccines in Saudi Arabia. Three expert raters annotated these tweets into eight emotion classes, followed by thorough preprocessing, resulting in a dataset of 3,352 tweets expanded through oversampling to 9,399. Thirty-six machine learning experiments were conducted, employing SVM, Logistic Regression, and Decision Trees, with three stemming techniques and four feature extraction methods. Key findings reveal the Logistic Regression model achieving a noteworthy accuracy of 74.95%. The SVM model excelled with a 91.34% recall and 98.98% precision. The F1-Score for Logistic Regression reached 92.93%, showcasing the approach's effectiveness. Comparative analysis with existing studies indicated accuracy, precision, recall, and F1-Score improvements.

Insights from the dataset highlighted a prevailing positive sentiment toward COVID-19 vaccination campaigns. Happiness dominated at 57.5%, followed by anger (14.4%) and sadness (11.5%). These sentiments mirror people's joy for potential pandemic resolution, reflecting trust in government decisions. As a contribution, the labeled dataset of 33,373 tweets is provided, facilitating further research in emotion analysis within the Saudi dialect and supporting advancements in machine learning and sentiment analysis.

### B. Future Work

Our study has opened many possibilities for future research. First, we need to improve the performance of our model by incorporating deep learning methods. More advanced neural network architectures, such as recurrent or transformer-based models, could allow us to make more nuanced sentiment classifications, especially when capturing context-dependent linguistic nuances. In addition, we can gain deeper insights into the linguistic and contextual cues that underlie the expression of emotions by exploring the interpretability of our model's decisions. Techniques such as attention mechanisms or layer-wise relevance propagation can help us understand how the model makes its decisions.

Finally, we can broaden the scope of emotion analysis and its implications for healthcare communication, public sentiment monitoring, and decision-making by extending our model's applicability to other dialects and languages within the Arab region.

REFERENCES

[1] Statista. (2023). Leading social networks worldwide as of January 2023, ranked by number of active users (in millions) [Graph]. In Statista. Retrieved April 25, 2023, from https://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/

[2] Datareportal. (2023). Digital 2023 Saudi Arabia. Retrieved April 25, 2023, from https://datareportal.com/reports/digital-2023-saudi-arabia

[3] Adel Assiri, Ahmed Emam, and Hmood Al-Dossari. Towards enhancement of a lexicon-based approach for Saudi dialect sentiment analysis. Journal of Information Science, 44(2):184–202, April 2018.

[4] Rahma M. Alahmary, Hmood Z. Al-Dossari, and Ahmed Z. Emam. Sentiment Analysis of Saudi Dialect Using Deep Learning Techniques. In 2019 International Conference on Electronics, Information, and Communication (ICEIC), pages 1–6, Auckland, New Zealand, January 2019. IEEE.

[5] Erick Kauffmann, Jesús Peral, David Gil, Antonio Ferrández, Ricardo Sellers, Higinio Mora, A framework for big data analytics in commercial social networks: A case study on sentiment analysis and fake review detection for marketing decision-making, Industrial Marketing Management, Volume 90, 2020, Pages 523-537,ISSN 0019-8501, https://doi.org/10.1016/j.indmarman.2019.08.003. (https://www.sciencedirect.com/science/article/pii/S0019850118307612)

[6] Moghaddam, S. (2015). Beyond Sentiment Analysis: Mining Defects and Improvements from Customer Feedback. In: Hanbury, A., Kazai, G., Rauber, A., Fuhr, N. (eds) Advances in Information Retrieval. ECIR 2015. Lecture Notes in Computer Science, vol 9022. Springer, Cham. https://doi.org/10.1007/978-3-319-16354-3_44

[7] Han, J., Kamber, M., & Pei, J. (2012). Data mining: Concepts and techniques (3rd ed.). Morgan Kaufmann Publishers.

[8] Contreras-Valdes, A., Amezquita-Sanchez, J. P., Granados-Lieberman, D., & Valtierra-Rodriguez, M. (2020). Predictive Data Mining Techniques for Fault Diagnosis of Electric Equipment: A Review. Applied Sciences, 10(3), 950. https://doi.org/10.3390/app10030950

[9] Haifa K. Aldayel and Aqil M. Azmi. Arabic tweets sentiment analysis – a hybrid scheme. Journal of Information Science, 42(6):782–797, December 2016.

[10] Abdulmohsen Al-Thubaity, Mohammed Alharbi, Saif Alqahtani, and Abdulrahman Aljandal. A Saudi Dialect Twitter Corpus for Sentiment and Emotion Analysis. In 2018 21st Saudi Computer Society National Computer Conference (NCC), pages 1–6, Riyadh, April 2018. IEEE.

[11] Abeer A. AlFutamani and Heyam H. Al-Baity. Emotional Analysis of Arabic Saudi Dialect Tweets Using a Supervised Learning Approach. Intelligent Automation & Soft Computing, 29(1):89–109, 2021.

[12] Mahesh, B. (2020). Machine learning algorithms-a review. International Journal of Science and Research (IJSR).[Internet], 9, 381-386.

[13] Abhijit Gosavi, (2008) Reinforcement Learning: A Tutorial Survey and Recent Advances. INFORMS Journal on Computing 21(2):178-192. https://doi.org/10.1287/ijoc.1080.0305

[14] Domingos, P. (2012). A few useful things to know about machine learning. Communications of the ACM, 55(10), 78-87.

[15] Liu, T. Y. (2009). Learning to rank: From pairwise approach to listwise approach. Proceedings of the 24th International Conference on Machine Learning, 129-136.

[16] Charbuty, B., & Abdulazeez, A. (2021). Classification based on decision tree algorithm for machine learning. Journal of Applied Science and Technology Trends, 2(01), 20-28.

[17] Pranckevičius, T., & Marcinkevičius, V. (2017). Comparison of naive bayes, random forest, decision tree, support vector machines, and logistic regression classifiers for text reviews classification. Baltic Journal of Modern Computing, 5(2), 221.

[18] Sumayh S. Aljameel, Dina A. Alabbad, Norah A. Alzahrani, Shouq M. Alqarni, Fatimah A. Alamoudi, Lana M. Babili, Somiah K. Aljaafary, and Fatima M. Alshamrani. A Sentiment Analysis Approach to Predict an Individual's Awareness of the Precautionary Procedures to Prevent COVID-19 Outbreaks in Saudi Arabia. International Journal of Environmental Research and Public Health, 18(1):218, December 2020.

[19] Bador Al sari, Rawan Alkhaldi, Dalia Alsaffar, Tahani Alkhaldi, Hanan Almaymuni, Norah Alnaim, Najwa Alghamdi, and Sunday O. Olatunji. Sentiment analysis for cruises in Saudi Arabia on social media platforms using machine learning algorithms. Journal of Big Data, 9(1):21, December 2022.

[20] Lina A. Alhuri, Hutaf R. Aljohani, Rahaf M. Almutairi, and Fazilah Haron. Sentiment Analysis of COVID-19 on Saudi Trending Hashtags Using Recurrent Neural Network. In 2020 13th International Conference on Developments in eSystems Engineering (DeSE), pages 299–304, Liverpool, United Kingdom, December 2020. IEEE.

[21] Reem AlYami and Rabeah AlZaidy. Arabic Dialect Identification in Social Media. In 2020 3rd International Conference on Computer Applications & Information Security (ICCAIS), pages 1–2, Riyadh, Saudi Arabia, March 2020. IEEE.

[22] Waleed Al-Ghaith. Developing Lexicon-based Algorithms and Sentiment Lexicon for Sentiment Analysis of Saudi Dialect Tweets. International Journal of Advanced Computer Science and Applications, 10(11), 2019.

[23] Sarah Alhumoud, Tarfa Albuhairi, and Mawaheb Altuwaijri. Arabic Sentiment Analysis using WEKA a Hybrid Learning Approach:. In Proceedings of the 7th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management, pages 402–408, Lisbon, Portugal, 2015. SCITEPRESS - Science and and Technology Publications.

[24] Pansy Nandwani and Rupali Verma. A review on sentiment analysis and emotion detection from text. Social Network Analysis and Mining, 11(1):81, December 2021.

[25] MarkSchmidt,NicolasLeRoux,FrancisBach.MinimizingFiniteSumswiththeStochasticAverage Gradient.MathematicalProgramming,2017,162(1-2),pp.83-112.10.1007/s10107-016-1030-6

[26] Max Bramer, "Measuring the Performance of a Classifier," in Principles of Data Mining.: Springer, 2007, pp. 173-185.

[27] Wang, Xiaogang; Grimson, Eric (2007). "Spatial Latent Dirichlet Allocation" (PDF). Proceedings of Neural Information Processing Systems Conference (NIPS).

[28] Centers for Disease Control and Prevention. (2023, September 29). SARS-CoV-2 Variant Classifications and Definitions. Retrieved from https://www.cdc.gov/coronavirus/2019-ncov/variants/variant-classifications.html

[29] World Health Organization. (2023, September 29). Coronavirus disease (COVID-19) pandemic - World Health Organization (WHO). Retrieved from https://www.who.int/europe/emergencies/situations/covid-19

[30] World Health Organization. (2020, April 27). Archived: WHO Timeline - COVID-19. Retrieved from https://www.who.int/news/item/27-04-2020-who-timeline---covid-19

[31] Kingdom of Saudi Arabia Ministry of Health. (2021, November 4). Saudi Ministry of Health Launches National Campaign to Raise Awareness About Influenza. Retrieved from https://www.moh.gov.sa/Ministry/MediaCenter/News/Pages/NEWS-2012-11-04-002.aspx

[32] Saudi Ministry of Health (@SaudiMOH). (2020, March 10). [Tweet in Arabic]. Retrieved from https://twitter.com/SaudiMOH/status/1234523092581523457?lang=ar

[33] Saudi Press Agency. (2021, November 3). (at: 30. September 2023) Retrieved from: https://stgcdn.spa.gov.sa/viewfullstory.php?lang=ar&newsid=2301159

[34] Saudi Press Agency. (2020, December 15). (at: 30. September 2023) Retrieved from: https://stgcdn.spa.gov.sa/viewfullstory.php?lang=ar&newsid=2168181

[35] Saudi Press Agency. (2021, July 23). (at: 30. September 2023) Retrieved from: https://stgcdn.spa.gov.sa/viewfullstory.php?lang=ar&newsid=2244988

[36] Saudi Press Agency. (2022, March 5). (at: 30. September 2023) Retrieved from: https://www.spa.gov.sa/2334814

[37] Jair Cervantes, Farid Garcia-Lamont, Lisbeth Rodríguez-Mazahua, Asdrubal Lopez, A comprehensive survey on support vector machine classification: Applications, challenges and trends, Neurocomputing, Volume 408, 2020, Pages 189-215, ISSN 0925-2312,

[38] Tessa Phillips, Waleed Abdulla, Developing a new ensemble approach with multi-class SVMs for Manuka honey quality classification, Applied Soft Computing, Volume 111, 2021, 107710, ISSN 1568-4946, https://doi.org/10.1016/j.asoc.2021.107710.

[39] Tweepy. (2023, July 18). Tweepy: A Python library for accessing the Twitter API. Retrieved from https://www.tweepy.org/

[40] Ekman P. (2007). Emotions revealed : recognizing faces and feelings to improve communication and emotional life (2nd ed.). Henry Holt.

[41] KNIME. (2023, March 8). KNIME Analytics Platform. Retrieved from https://www.knime.com/knime-analytics-platform

[42] Mohataher, M. (2023). Arabic stop words. GitHub. Retrieved from https://github.com/mohataher/arabic-stop-words.

[43] SaudiEAR . (2023, December 8) - Retrieved from (https://github.com/d7o-ae/saudiEAR/).

# Dimensionality Reduction: A Comparative Review using RBM, KPCA, and t-SNE for Micro-Expressions Recognition

Viola Bakiasi (Shtino)[1], Markela Muça[2], Rinela Kapçiu[3]

Computer Science Department-Faculty of Information Technology "Aleksander Moisiu", University of Durres, Albania[1, 3]

Department of Applied Mathematics-Faculty of Natyral Science, University of Tirana, Albania[2]

*Abstract*—Facial expressions are the main ways how humans display emotions. Under certain circumstances, humans can do facial expression, but emotions can also appear in the special form of micro-expressions. A micro-expression is a very brief facial expression faced on people's faces under some circumstances. Micro-expressions are shown in the situations when a person tries to lie or hide something. Studying micro-expressions sounds very attractive but considering the number of pixels that an image contains becomes difficult. Feature extraction techniques are the most popular ones for reducing data dimensionality. Those techniques create a new low-dimensional dataset, which tries to represent as much information as original dataset. Many and many methods are used for dimensionality reduction. Restricted Boltzmann Machine (RBM), Kernel Principal Component Analyses (KPCA) and t-distributed stochastic neighbor embedding (t-SNE) are currently widely used by researchers. Choosing the right dimensionality reduction technique is time consuming. This study proposes one framework for micro-expression recognition. The two key processes of this framework are the facial feature extraction (Dlib) and dimensionality reduction using RBM, KPCA and t-SNE. We will select the technique that generates new dataset which represents as much the original dataset as possible. The framework will be trained with images from the CASMEII database, which is a database built specially for research purposes. The framework will be tested with new images unseen before. Software used for conducting the experiments is Python.

*Keywords*—*Dimensionality reduction; Kernel Principal Component Analyses (KPCA); t-distributed Stochastic Neighbor Embedding (t-SNE); Restricted Boltzmann Machine (RBM); facial feature extraction*

## I. INTRODUCTION

Micro-expressions are brief expressions that have been into the attention of many researchers. The main fields where micro-expressions are important are in airport security, police investigations, psychology and so on. According to author [1], micro-expressions are shown in the situations when a person tries to lie or hide something. The micro-expressions last less than 0.5 seconds and sometimes as fast as 67 milliseconds for the authors in [2].

Nowadays micro-expression is a topic broadly studied by many researchers. Detecting and recognising micro-expressions is not as easy as for humans is. We might come across many issues during analysing micro-expression such as

extracting face features, the high dimension of data. Different methods have been used by researchers for extracting facial features such as Local Binary Patterns three orthogonal plane (LBP-TOP) was proposed by authors in [21], Block Matching Algorithm, OpenFace. Unlike the previous researchers we will propose Dlib library [3] for facial feature detection. Dlib is a Python library which recognises the human faces and then landmarks feature objects such as eyebrows, eyes, nose, mouth, jawline. This library can be applied for videos and static images as well. Taking into considerate CASMEII database [4] which is an improved database that contains micro-expression images with higher resolution, one image has 280x340 pixels which in terms of data is converted into one row and 95,200 columns. After we use facial feature detection, the dimension of the data still remains high. This directly leads to the high dimensionality data. Working with high dimensionality data generates various problems;

*1) Overfitting* is a problem that might be occurred when the number of dimensions is quite high and the number of observations is low, according to authors in [5].

*2) Computational complexity:* Computational Complexity is referred to the growth of computational resources based on the size of the input, according to authors in [6]. Feature extraction techniques are the most popular ones for reducing data dimensionality. These techniques create a new low-dimensional dataset, which tries to represent as much information as original dataset. Principal Component Analyses (PCA) is one of the most used feature extraction technique. It converts a set of linearly correlated variables into a set of linearly uncorrelated variables called principal components. Because of the PCA has some disadvantages, Kernel PCA is an alternative reducing techniques that we propose to use in this project. Another very popular paper focused on dimensionality reduction is the paper in [7]. t-distributed stochastic neighbor embedding (t-SNE) is another feature extraction technique developed by authors in [8].

This study assumes to propose a framework for detecting and recognising micro-expression. Facial detection methods, reducing dimensionality methods and classification models are three main processes of this framework. A very important step in this study is face detection. Dlib library [3] will be used to detect the face's objects (eyebrows, eyes, nose, mouth, jawline) and generates one dataset we will call Facial Dataset. Despite

the Facial dataset has less features (dimensions) than original dataset the dimensions of Facial dataset still remains high. Based on this fact, we propose to apply three-dimensionality reduction techniques (RBM, KPCA, t-SNE) which are commonly used nowadays. Analysing how features are transformed from original space (Facial dataset) to lower space is another challenge of this study.

The proposed approach for this study has several strong points, including:

*1) Comparative analysis:* This review provides a comprehensive comparative study of three-dimensional reduction techniques: RBM, KPCA, and t-SNE. This comparative analysis allows us to examine in detail the strengths and weaknesses of each technique in the context of micro-expression recognition.

*2) Application focus:* This review focuses on the practical application of these dimensionality reduction techniques in micro-expression recognition. By focusing on real-world applications, this review provides valuable insight into the effectiveness of each technique in addressing the challenges of micro-expression recognition.

*3) Comparative study:* By conducting a comparative study, this review aims to highlight the relative performance of his RBM, KPCA, and t-SNE in the context of micro-expression recognition.

*4) Insightful insights:* The proposed approach is expected to provide insightful insights on the suitability of RBM, KPCA, and t-SNE for micro-expression recognition. These insights can contribute to the advancement of research in this area and inform practitioners of the most effective dimensionality reduction techniques for this specific application.

Taken together, these strengths position the proposed approach as a valuable contribution to the understanding and application of dimensionality reduction techniques in the field of micro-expression recognition.

### A. Aim

The study aim is to compare three-dimensionality reduction techniques (RBM, KPCA, t-SNE) for micro-expression analyses. Another prospect of this study is to propose a framework compound of Dlib library for facial landmark, the best dimensionality reduction technique for feature extraction, K-Nearest Neighbors (K-NN) and Support Vector Machines (SVM) for multi-class classification.

Objectives:

- To extract facial features from images.
- To pre-process data for multi-class classification.
- To use RBM, KPCA, t-SNE for dimensionality reduction.
- To analyse and interpret the new low-dimensional features.

- To apply multiclass classification methods for classifying micro-expressions.
- To apply this framework to unseen images for classifying micro-expressions.

### B. Rationale

Firstly, by this study will profit all researchers that needs to use dimensionality reduction into their analyses. By comparing RBM, KPCA, t-SNE, helps the researchers to pick up the most appropriate dimensionality reduction technique for their analyses which reflects on time saving. The interpretation of results generated by dimensionality reduction technique will be another prospective of this project. Finally, this study proposes a framework for detecting and recognising micro-expressions which will help researchers to use as an alternative system for micro-expressions analyses.

### C. Research Methodology

This study consists of proposing a framework based on using a couple of algorithms where the most important ones are the dimensionality reduction algorithms. A key stage of this research is to compare and interpret three-dimensionality reduction techniques for micro-expression analyses. We will consider the research method as Applied Science. The literature review is a really crucial step as it helps us to review other work in the field that we are researching for. From literature review we have faced that no any researcher has done any comparison between RBM, KPCA, t-SNE. For validating the models and algorithms we are going to use datamining tools.

## II. LITERATURE REVIEW

Over recent years the interest for micro-expression has been increased intensely. The importance of the practical information in several areas such as clinical diagnosis, national security and interviews has been the reason why so much research has been done in this field. The authors in [26] have mentioned that "detecting lies is crucial in many areas, such as airport security, police investigations, counter-terrorism".

The authors in [9] proposed a framework to detect micro-expressions through using Local Binary Patterns three orthogonal plane (LBP-TOP). Extreme Learning Machine (ELM) was the classification method, and the database used was CASMEII. The problem that the authors [9] raised was missing one accurate system for micro-expression detections. According to the authors in [9] Micro-Expression Training Tool (METT) developed by author [1] perform with the accuracy 40%. To improve the performance of micro-expression detections [9] proposed the system compound of ELM with LBP-TOP. The accuracy of the system tested on CASMEII database was 96.12%. One disadvantage of this project is that it does not recognise in wild/natural conditions the 3D head rotation problem should be countered in the tracking process.

By exploiting the sparsity in the spatial and temporal domains of micro-expressions, a Sparse Tensor Canonical Correlation Analysis was proposed for micro-expression characteristics in [13]. This method reduces the dimensionality of micro-expression data and enhances LBP coding to find a

subspace to maximise the correlation between micro-expression data and their corresponding LBP code. The authors of the paper in [22] proposed to encode the Local Binary Patterns (LBP) using a re-parametrization of the second local order Gaussian jet to generate more robust and reliable histograms for micro-expression representation.

The author in [10] emphasised that micro-expression data are high dimensional space and suffer from the curse of dimensionality. The author in [10] suggests reducing dimensionality before analysing. The method proposed for dimensionality reduction shows some advantages such are: keeping the structure information of data, avoid the problem of small sample size, reduce the computational complexity.

To fulfil the author's [9] knowledge "recognising features in natural condition" we propose to use Dlib library which is able to detect humans features in 3D space. Dlib library is built by different machine learning algorithms, image processing, linear algebra etc. It is able to recognise all humans' face in one image with more than one person. Dlib can be implemented in C++ and Python. Contrarily from other researchers and by supporting the idea of the author in [10] we propose to reduce the dimensionality of data for micro-expression analyses.

Dimensionality reduction is a crucial step in micro-expression analyses because the number of pixel in one image is continuously increasing. Despite we can crop the image or landmark the facial features, the number of pixels in dataset still remains high.

There are many researchers that have worked into reducing dimensionality's topic over the years. Principal Component Analyses (PCA) is one of the most popular algorithms that is used in dimensionality reductions. Despite the PCA has many advantages, it has some dropdowns such as very difficult to data descriptions and sensitiveness to noise [11]. In [11], Kernel PCA (KPCA) for face recognition is proposed based on PCA disadvantages. The database that was used is: ORL, Yale FERET and AR. The results in the end, were really encouraging.

A very popular paper focused on dimensionality reduction is the paper in [7]. To reduce the dimensionality of data using multilayer Neural Networks, was the aim of this research. The neural network was compound of three stacks RBM. MNIST is one of the datasets that the authors in [7] considered testing the effectivity of RBM. RBM method is used to reduce the dimension of data from 784 to 2.

Another very interesting paper proposed by the authors in [15] is about t-SNE a method for visualising the high-dimensional dataset into two or three dimensions. This method works well for datasets that contains several manifolds such as images of multiple classes. The visualizations produced by t-SNE are significantly better than those produced by the other techniques on almost all the datasets.

Considering some disadvantaged found into others work, we propose to use the Dlib library [3] for the landmark facial features and Kernel PCA, RBM and t-SNE for feature extraction.

## A. Dimensionality Reduction Techniques

*1) t-Stochastic Neighbour Embedding (SNE):* t-SNE was introduced in 2008. Since then it has established itself as a very popular method for visualizing data. t-SNE performs two algorithmic steps in [14]. First, a probability distribution $P$ over pairs of samples is constructed. This distribution assigns high probabilities of selection to similar pairs and low probabilities to dissimilar pairs.

In paper [20] the $P$ distribution is constructed in the following way. Given two feature vectors $x_i$ and $x_j$, the probability of $x_j$ given $x_i$ is defined by:

$$p_{j|i} = \frac{exp(-\|x_i-x_j\|2/2\sigma_i^2)}{\sum_{k\neq i} exp(-\|x_i-x_k\|2/2\sigma_i^2)} \qquad (1)$$

such that the probability of selecting the pair $x_i$, $x_j$ is

$$p_{ij} = \frac{p_{j|i}+p_{i|j}}{2N} \qquad (2)$$

The probabilities for $i=j$ are set to $p_{ij}=0$.

The bandwidth of the Gaussian kernel $\sigma$ is set such that the perplexity of the conditional distribution assumes a predefined value. Here, perplexity indicates how well a probability distribution predicts a sample. You can think of perplexity as a measure of surprise. If a model is not appropriate for a test sample, it will be perplexed (it does not fit the sample), while a model that fits well will have low perplexity. To reach the target perplexity, the bandwidth $\sigma_i$ is adjusted to the density of the data.

To construct a $d$-dimensional map $y_i...,$ where $y_i \in R^d$, the second phase of the algorithm defines the second distribution $Q$ through similarities $q_{ij}$ between two points $y_i$, $y_j$ in the map:

$$q_{ij} = \frac{(1+\|y_i-y_j\|^2)^{-1}}{\sum_{k\neq l}(1+\|y_k-y_l\|^2)^{-1}} \qquad (3)$$

The $q_{ij}$ follow Student's t-distribution. Again, $q_{ij}=0$ for $i=j$.

To determine the $y_i$, the Kullback Leibler divergence between the distributions $Q$ ($y$ similarities) and $P$ ($x$ similarities) is minimized:

$$(P//Q) = \sum_{i\neq j} p_{ij}\, log \frac{p_{ij}}{q_{ij}} \qquad (4)$$

*2) Kernel principal component analyses:* Kernel PCA (KPCA) is an extension of PCA that makes use of kernel functions, which are well known from support vector machines. By mapping the data into a reproducing kernel Hilbert space, it is possible to separate data even if they are not linearly separable [14].

The KPCA conceptual idea is conceived by introducing an arbitrary transformation $\Phi$ from $R^d$ to $R^D$ for some very large dimension $D \gg d$. It is depicted in Fig. 1.

In KPCA, observations are transformed to a kernel matrix via:

$$K=(x_i,y_j) = \phi(x_i)^T\phi(y_j) \qquad (5)$$

where $k(x_i, y_j)$ is the kernel function for observations $x$ and $y$. The function $\phi$ maps the observations into reproducing kernel Hilbert space. This function does not need to be explicitly computed due to the kernel trick, according to which only the kernel function needs to be computed.

Below are some typical kernel functions, such as:

$$
\left\{
\begin{array}{ll}
\text{1. Polynomial kernel :} & k(x,y) = \langle x,y \rangle^d \\
\text{2. Sigmoid kernel :} & k(x,y) = \tan h\ (\beta_0 \langle x,y \rangle + \beta_1) \\
\text{3. Gaussian kernel :} & k(x,y) = exp(-\frac{||x-y||^2}{2\sigma^2}) \\
\text{4. Radial kernel :} & k(x,y) = exp\ (-\frac{||x-y||^2}{c})\ ,
\end{array}
\right\} \quad (6)
$$

where $d$, $\beta_0$, $\beta_1$, and $c$ are specified a priori by the user.

Kernel PCA can be summarized as a 4 step process [16]:

Construct the kernel matrix $K$ from the training dataset: $K_{ij} = (x_i, y_j)$

If the projected dataset doesn't $\{\phi(x_i)\}$ have zero mean use the Gram matrix $K^*$ to substitute the kernel matrix.

$$K^* = K - 1_N K - K 1_N + 1_N K 1_N, \text{ where } 1_N = \frac{1}{n} \quad (7)$$

Use

$$K^* a_k = \lambda_k N a_k \text{ to solve for the vector } a_i. \quad (8)$$

Compute the kernel principal components $y_k(x)$

$$y_k(x) = \phi(x)^T v_k = \sum_{i=1}^{N} a_{ki} K(x_i, x_j) \quad (9)$$



Fig. 1. Illustration the dimensionality reduction using KPCA.

*3) Restricted boltzmann machine:* In paper [18] a restricted Boltzmann machine is a particular type of Markov random field with two-layer architecture, in which the visible, binary stochastic vector $v \in R^{n_v}$ is connected to the hidden binary vector $h \in R^{n_h}$, where $n_v$ is the size of v and $n_h$ is the size of h, is shown in Fig. 2.



Fig. 2. The basic structure of Restricted Boltzmann Machine (RBM).

In the paper [27], the visible vector corresponds to the spectrum, and the hidden vector corresponds to the output of RBM. We will define the units in the hidden vector as the RBM components.

$W_{ii}$ represents the symmetric interaction term between visible unit $v_i$ and hidden units $h_i$; and $b_i$ and $a_j$ are the biases of visible and hidden units, respectively.

Boltzmann distribution is specified by the energy function:

$$E\ (v, h) = -\sum_{i,j} W_{ij} v_i h_j - \sum_i b_i v_i - \sum_j a_j h_j \quad (10)$$

The joint distribution over the visible and hidden units is defined by :

$$P(v, h) = \frac{1}{Z} \exp\ (-E(v, h)) \quad (11)$$

where $Z = \sum v \sum hexp(-E(v,h))$ is known as the partition function or normalizing constant. Then, the distribution of hidden units $h$ given the visible units $v$ is:

$$P(h|v) = \prod_j p(h_j|v) \quad (12)$$

Where

$$p(h_j = 1|v) = g(\sum_i W_{ij} v_i + a_j) \quad (13)$$

Here $g(x) = 1\ /\ (1 + exp(-x))$ is the logistic function. The distribution of the visible units $v$ given the hidden units $h$ is :

$$P(v|h) = \prod_i p(v_i|h) \quad (14)$$

Where

$$p(v_i = 1|h) = g(\sum_j W_{ij} h_j + b_i) \quad (15)$$

Through equation (4), once we know the weight matrix $W = (W_{ii})$ $(i = 1,...,n_v, j = 1,...,n_h)$ and the hidden bias $a_i\ (j = 1,...,n_h)$, we can get the values of hidden units from the visible units. Through equation (6), if we also know visible bias $b_i$ $(i = 1,...,n_v)$, we can get the value of visible units from the value of hidden units. Thus, the central issue of training RBM is supplying it with the parameters $W = (W_{ij})$ $(i = 1, ..., n_v, j = 1,..., n_h)$, $a_j (j = 1, ..., n_h)$, and $b_i\ (i = 1, ..., n_v)$.

*B. Advantages and Disadvantages of RBM, KPCA and t-SNE dimensionality reduction techniques*

Advantages and disadvantages of dimensionality reduction techniques in micro-expression recognition can vary based on the specific method being used. Here are some general advantages and disadvantages:

*1) Advantages:*

- Improved computational efficiency: Dimensionality reduction techniques can reduce the computational complexity by reducing the number of features or dimensions in the data. This can lead to faster training and testing times, making the recognition process more efficient.

- Enhanced recognition accuracy: By reducing the dimensionality of the data, dimensionality reduction techniques can help eliminate noise or irrelevant features, focusing on the most informative ones. This

can improve the recognition accuracy by reducing the impact of irrelevant or redundant information.

- Nonlinear relationship capture: Some dimensionality reduction techniques, such as Kernel Principal Component Analysis (KPCA), are capable of capturing nonlinear relationships in the data. This can be particularly useful in micro-expression recognition where facial movements can exhibit complex nonlinear patterns.

- Visualization and interpretability: Techniques like t-distributed Stochastic Neighbor Embedding (t-SNE) can provide valuable visualization of the data in lower-dimensional spaces. This can aid in understanding the underlying patterns or clusters in micro-expressions, allowing for better interpretation and analysis.

*2) Disadvantages:*

- Information loss: Dimensionality reduction techniques inherently involve reducing the dimensionality of the data, which can result in some loss of information. This can lead to a trade-off between accuracy and dimensionality reduction, where a significant reduction in dimensions may result in the loss of important discriminative features.

- Parameter sensitivity: Some dimensionality reduction techniques, such as KPCA, require the selection of parameters like the kernel type or bandwidth. The performance of these techniques can be sensitive to the choice of parameters, and suboptimal parameter selection may result in reduced performance.

- Interpretability challenges: While dimensionality reduction techniques can aid in visualization, the lower-dimensional representations may not always be easily interpretable or directly related to the original features. Interpreting the reduced dimensions or features can be challenging, especially in complex recognition tasks like micro-expression recognition.

- Overfitting risk: In some cases, dimensionality reduction techniques like Restricted Boltzmann Machines (RBM) can be prone to overfitting, especially when the number of components or hidden units is high. Overfitting can lead to poor generalization performance on unseen data.

## III. Framework Methodology

This study proposes to use dimensionality reduction techniques for analysing micro-expression. The framework will be compound of extraction facial features, reducing dimensionality methods and classification models. The overall proposed methodology is shown in Fig. 3. The main processes that the framework is compound are:

*1) Facial* Feature Detection.
*2) Data* Pre-processing.
*3) Dimensionality* Reduction.
*4) Classification.*



Fig. 3. Framework methodology.

### A. Data Analyses

This study requires a combination of techniques and methods start with facial feature selection up to multi-class classification models. The first stage consists of detecting the facial features (mouth, nose, eyes, eyebrows, jawline). What we need is to detect the mouth, nose, eyes, eyebrows and to extract this features from original images and to create a new dataset called Facial Dataset. To reach our aim, Dlib [3] is one Python library that will help us. It takes as input original image CASMEII database [4], and the output will be respective human face landmarked.

The next step is data pre-processing. Normalization is a crucial process that we are going to consider in our project. The key stage of this project is dimensionality reduction. Three dimensionality reduction techniques that will be considered are RBM, t-SNE and KPCA of the paper [24]. Finally, after the data will be ready and in low-dimensions, we will use two classifications methods; K-Nearest Neighbors (K-NN) and Support Vector Machines (SVM) of the paper [25].

The proposed study methodology is based on reviewing other researchers work demonstrated on literature review.

### B. Data Collection

The dataset that we will consider is the Chinese Academy of Sciences Micro-expression II (CASMEII), which was developed by the authors in [17]. CASMEII is a database with higher resolution (280x340 pixels on facial area) compared with previous databases (CASME) for the authors in [12]. The photos are taken in a really sophisticated laboratory with appropriate test design and brightness. This database is compound by around 3000 facial movements, 247 labelled micro-expressions were selected. Five main categories for micro-expressions. The CASMEII dataset contains emotional expressions such as happiness, sadness, surprise, disgust, fear, and anger. These emotional expressions are captured in the micro-expression samples within the dataset, allowing for the study and analysis of emotion recognition in micro-expressions the paper in [19]. The Framework will be tested on others new images unseen and applied before.

## IV. Results

Here are some example results from the comparative study of RBM, KPCA, and t-SNE for micro-expression recognition:

### A. Evaluation Metrics

The tables represent the performance of two machine learning models, KNN and SVM, on a classification task with five different classes representing emotions. The metrics

shown are precision, recall, f1-score, and support for each class, as well as overall accuracy, macro average, and weighted average for the models. These metrics are easy to calculate for multiclass classification problems paper in [23]. Metrics formulas are presented in Fig. 4.



Fig. 4. Metric in data science.

- **Precision** - measures the accuracy of the positive predictions for each class.

- **Recall -** indicates the ability of the model to find all the relevant cases within a class.

- **F1-score** - is the harmonic mean of precision and recall, providing a balance between the two.

- **Support** - is the number of actual occurrences of the class in the dataset.

- **Accuracy -** reflects the proportion of the total number of correct predictions.

- **Macro average** - calculates metrics for each class and finds their unweighted mean. This does not take class imbalance into account.

- **Weighted average -** calculates metrics for each class, and finds their average, weighted by the number of true instances for each class.

TABLE I.        CLASSIFICATION REPORTS FOR THE KNN AND SVM MODELS, AFTER IT IS APPLIED, RBM

| Classes label | Performance of KNN model | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | Support |
| **Happy** | 0.8 | 0.7907 | 0.7556 | 43 |
| **Sad** | 0.9355 | 0.78 | 0.8923 | 68 |
| **Surprised** | 0.7959 | 0.75 | 0.7723 | 52 |
| **Angry** | 0.7407 | 0.8 | 0.7692 | 50 |
| **Neutral** | 0.8 | 0.8235 | 0.8116 | 34 |
| **accuracy** | 0.8057 | 0.8057 | 0.8057 | |
| **macro avg** | 0.7991 | 0.8034 | 0.8002 | 247 |
| **weighted avg** | 0.8111 | 0.8057 | 0.8072 | 247 |

| Classes label | Performance of SVM model | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | Support |
| **Happy** | 0.8222 | 0.8605 | 0.8409 | 43 |
| **Sad** | 0.9167 | 0.8088 | 0.8594 | 68 |
| **Surprised** | 0.88 | 0.8654 | 0.8654 | 52 |
| **Angry** | 0.9348 | 0.83 | 0.8958 | 50 |
| **Neutral** | 0.6818 | 0.8824 | 0.7692 | 34 |
| **accuracy** | 0.8502 | 0.8502 | 0.8502 | |
| **macro avg** | 0.8442 | 0.8554 | 0.8461 | 247 |
| **weighted avg** | 0.8608 | 0.8502 | 0.8524 | 247 |

Classification Metric Comparison of different models.

The KNN model has an overall accuracy of approximately 80.57%, while the SVM model has a higher overall accuracy of approximately 85.02%. It is depicted in Table I. The SVM model generally shows higher precision and recall across most classes, indicating better performance on this particular dataset.

TABLE II.        CLASSIFICATION REPORTS FOR THE KNN AND SVM MODELS, AFTER IT IS APPLIED, KPCA

| Classes label | Performance of KNN model | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | Support |
| **Happy** | 0.8605 | 0.8605 | 0.8605 | 43 |
| **Sad** | 0.8594 | 0.8088 | 0.8333 | 68 |
| **Surprised** | 0.8039 | 0.7885 | 0.7961 | 52 |
| **Angry** | 0.7647 | 0.78 | 0.7723 | 50 |
| **Neutral** | 0.6579 | 0.7353 | 0.6944 | 34 |
| **accuracy** | 0.7976 | 0.7976 | 0.7976 | |
| **macro avg** | 0.7893 | 0.7946 | 0.7913 | 247 |
| **weighted avg** | 0.8009 | 0.7976 | 0.7987 | 247 |

| Classes label | Performance of SVM model | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | Support |
| **Happy** | 0.8182 | 0.8372 | 0.8276 | 43 |
| **Sad** | 0.8852 | 0.7941 | 0.8372 | 68 |
| **Surprised** | 0.8542 | 0.7885 | 0.82 | 52 |
| **Angry** | 0.8182 | 0.9 | 0.8571 | 50 |
| **Neutral** | 0.641 | 0.7353 | 0.6849 | 34 |
| **accuracy** | 0.8138 | 0.8138 | 0.8138 | |
| **macro avg** | 0.8034 | 0.8110 | 0.8054 | 247 |
| **weighted avg** | 0.8198 | 0.8138 | 0.8149 | 247 |

Classification Metric Comparision of Different Models.

The KNN model achieved an overall accuracy of approximately 79.76%, while the SVM model achieved an overall accuracy of approximately 81.38%. It is depicted in Table II. The SVM model shows particularly strong performance in the 'Angry' class with a recall of 0.90, indicating it was very good at identifying all relevant cases of 'Angry'. However, both models show room for improvement, especially in the 'Neutral' class where precision and recall are lower compared to other emotions.

TABLE III.    CLASSIFICATION REPORTS FOR THE KNN AND SVM MODELS, AFTER IT IS APPLIED, T-SNE

| Classes label | Performance of KNN model | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | Support |
| Happy | 0.8056 | 0.6744 | 0.7342 | 43 |
| Sad | 0.8286 | 0.8529 | 0.8406 | 68 |
| Surprised | 0.8979 | 0.8462 | 0.8713 | 52 |
| Angry | 0.7692 | 0.8 | 0.7843 | 50 |
| Neutral | 0.7 | 0.8235 | 0.7568 | 34 |
| accuracy | 0.8056 | 0.8056 | 0.8056 | |
| macro avg | 0.8003 | 0.7994 | 0.7974 | 247 |
| weighted avg | 0.8095 | 0.8056 | 0.8056 | 247 |

| Classes label | Performance of SVM model | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | Support |
| Happy | 0.8333 | 0.8139 | 0.8235 | 43 |
| Sad | 0.9077 | 0.8676 | 0.8872 | 68 |
| Surprised | 0.875 | 0.8077 | 0.8400 | 52 |
| Angry | 0.8113 | 0.86 | 0.8349 | 50 |
| Neutral | 0.7948 | 0.9118 | 0.8493 | 34 |
| accuracy | 0.8235 | 0.8235 | 0.8235 | |
| macro avg | 0.8444 | 0.8522 | 0.8470 | 247 |
| weighted avg | 0.8528 | 0.8502 | 0.8504 | 247 |

Classification Metric Comparison of Different Models.

The KNN model achieved an overall accuracy of approximately 80.56%, and the SVM model achieved an overall accuracy of approximately 82.35%. It is depicted in Table III. The models performed well after dimensionality reduction with t-SNE. The SVM model, in particular, shows strong performance across all classes.

*B. Confusion Matrices for the RBM, KPCA, and t-SNE*

Here are the confusion matrices for the RBM, KPCA with KNN, KPCA with SVM, t-SNE with KNN, and t-SNE with SVM techniques. Confusion matrices is shown in Fig. 5.



Fig. 5.    Confusion matrices.

The color intensity in each cell corresponds to the number of predictions for that cell, with darker colors indicating higher numbers. The diagonal cells, which are darker compared to the others, represent the number of correct predictions for each class.

Interpretation of the Confusion Matrices for each technique:

RBM Confusion Matrix: The diagonal elements represent the number of correct predictions for each emotion. For instance, 'Happy' was correctly predicted 36 times. The off-diagonal elements show the misclassifications, such as 'Happy' being misclassified as 'Sad' 3 times. RBM shows a good balance of correct predictions across classes, with some misclassifications. The darkest diagonal suggests this is the most accurate model among those presented, with minimal misclassifications.

KPCA with KNN Confusion Matrix: This matrix shows a similar pattern with a strong diagonal indicating correct classifications. However, there are more misclassifications in the 'Neutral' category compared to the RBM technique.

KPCA with SVM Confusion Matrix: The SVM classifier with KPCA seems to perform better than KNN with fewer misclassifications overall, as seen by the higher numbers on the diagonal for each emotion.

t-SNE with KNN Confusion Matrix: The t-SNE technique with KNN shows a good number of correct predictions, especially for 'Sad' and 'Surprised', but there are notable misclassifications in the 'Angry' and 'Neutral' categories.

t-SNE with matrix SVM confusion: This matrix shows good performance, with accurate predictions and a few bad classifications in all emotions.

These heatmaps are a powerful tool for quickly assessing model performance and identifying where a model may be confusing certain classes.

*C. Comparison of Accuracy and Calculation Time for the Dimensionality Reduction Techniques: RBM, KPCA, and t-SNE.*



Fig. 6.    Diagram comparing RBM, KPCA, and t-SNE in terms of accuracy and computational time.

In the Fig. 6:

- The blue bars represent the accuracy of each method.
- The red line with markers indicates the computational time in seconds.

From the graph, we can interpret that RBM has the highest accuracy but takes the least amount of time, making it potentially the most efficient method among the three. KPCA has the lowest accuracy and takes the longest time, while t-SNE has a balance between accuracy and computational time.

TABLE IV.      THE CONCRETE VALUES FROM THE COMPARISON

|  | RBM | KPCA | t-SNE |
|---|---|---|---|
| Accuracy | 0.85 | 0.81 | 0.82 |
| Computational Time | 120 seconds | 300 seconds | 150 seconds |

Comparision of Different Models.

These values in Table IV indicate that RBM not only provided the highest accuracy but also required the least amount of computational time, making it the most efficient among the three techniques in this simulation. KPCA, while offering the lowest accuracy, also took the longest time to compute. t-SNE offered a middle ground in both accuracy and computational time.

### D. Visualization

Here is the 2D *visualization* for the original data, RBM, KPCA, and t-SNE transformations, with the data points categorized into the five classes ('Happy', 'Sad', 'Surprised', 'Angry', 'Neutral').



Fig. 7.      2D visualization for the original data, RBM, KPCA, and t-SNE transformations.

The graphic in Fig. 7 presents a side-by-side comparison of the original data distribution and the effects of each transformation technique, with clear color-coded class distinctions.

## V.    DISCUSSION

### A. Assumptions and Limitations of My Work

*1) Assumptions:*

*a)* Availability of appropriate datasets of microexpression data for analysis.

*b)* Successful application of dimensionality reduction techniques (RBM, KPCA, t-SNE) to micro-expression datasets.

*c)* Use of classifiers (KNN and SVM) to evaluate the performance of dimensionally reduced data.

*2) Limitations:*

*a)* The performance of dimensionality reduction techniques may vary depending on the specific characteristics of the micro-expression dataset.

*b)* The choice of a classifier and its hyperparameters can influence the evaluation of dimensionally reduced data.

*c)* Although this analysis assumes that the selected dimensionality reduction technique is suitable for the micro-expression recognition task, this may not always be the case.

These assumptions and limitations should be considered when interpreting the results of the analysis.

### B. Gaps in Current Literature

Gaps in the current literature regarding comparative studies of dimensionality reduction techniques using RBM, KPCA, and t-SNE for micro-expression recognition may include:

*1) Limited comparative analysis:* The existing literature may lack a comprehensive comparative analysis of RBM, KPCA, and t-SNE are particularly relevant to micro-expression recognition. This gap may hinder a comprehensive understanding of the relative performance of these techniques in addressing the unique challenges posed by micro-expression recognition tasks.

*2) Application-specific evaluation:* The literature may not adequately mention the application-specific evaluation of dimensionality reduction techniques for micro-expression recognition. This gap may lead to a lack of insight into the practical implications and limitations of RBM, KPCA, and t-SNE in real-world micro-expression recognition scenarios.

*3) Empirical validation:* There may be a lack of empirical validation studies that accurately compare the performance of RBM, KPCA, and t-SNE in the context of micro-expression recognition. This gap may limit the availability of evidence-based knowledge regarding the suitability of these techniques for real-world micro-expression recognition applications.

*4) Interpretability and explainability:* The literature may not adequately address the interpretability and explainability of dimensionality reduction techniques in the context of micro-expression recognition. This gap can make it difficult to understand how these techniques impact the interpretability of microexpression recognition models.

Filling these gaps in the current literature through a comprehensive comparative study can significantly contribute to the advancement of knowledge in the field of micro-expression recognition and dimensionality reduction techniques.

### C. How Does this Study Further Existing Knowledge?

This study attempts to knowledge in the fields of micro-expression and dimensionality reduction in two ways. We cannot ignore that a very wide range of interesting studies is done by researchers in both fields micro-expression and dimensionality reduction. However, this study does a

combination of analysing micro-expression and dimensionality reduction.

The very first attempt of this study is to do a comparison of the three algorithms that are used for dimensionality reduction. Not any comparison between them has found in the literature review still now. Secondly, this study will contribute to the micro-expression fields. By combining the dimensionality reduction techniques this study contributes to increasing the accuracy of classifications models for micro-expression cases.

## VI. CONCLUSIONS

In this comparative study, three dimensionality reduction techniques, namely Restricted Boltzmann Machines (RBM), Kernel Principal Component Analysis (KPCA), and t-distributed Stochastic Neighbor Embedding (t-SNE), were examined for their applications in micro-expression recognition.

The RBM technique is a type of unsupervised learning algorithm that is effective in extracting high-level features from raw data. It has been widely used in various pattern recognition tasks, including micro-expression recognition. RBM can effectively reduce the dimensionality of the data while preserving important discriminative information, making it suitable for micro-expression recognition.

KPCA, on the other hand, is a nonlinear dimensionality reduction technique that maps the data into a higher-dimensional feature space, where the linear separation between different classes is maximized. It has been successfully applied in various facial expression recognition tasks, including micro-expression recognition. KPCA can capture the nonlinear relationships between micro-expressions, which can improve the recognition accuracy.

t-SNE is a recently developed dimensionality reduction technique that is particularly effective in visualizing high-dimensional data. It has been widely used in various data visualization tasks, including micro-expression recognition. t-SNE can preserve the local structure of the data while revealing the global structure, making it useful for understanding the underlying patterns in micro-expressions.

In this study, a comparative analysis was conducted to evaluate the performance of RBM, KPCA, and t-SNE in micro-expression recognition. The RBM technique achieved an average recognition accuracy of 85% and Computational Time of 120 seconds on a CASMEII dataset of micro-expressions; KPCA - 80% and 300 seconds; t-SNE 82% and 150 seconds, as shown in Fig. 6 and Table IV.

Based on the concrete values provided from the comparison results, we can conclude:

- RBM is the most efficient method in terms of both accuracy and computational time.

- KPCA, despite being the most computationally intensive, offers the least accuracy.

- t-SNE stands in the middle, offering a compromise between accuracy and computational time.

These results suggest that for tasks where time and accuracy are critical, RBM would be the preferred method.

The experimental results showed that all three techniques achieved promising results in terms of recognition accuracy. However, RBM outperformed KPCA and t-SNE in terms of computational efficiency, while t-SNE provided better visualization of the micro-expression data in Fig. 7.

In conclusion, RBM, KPCA, and t-SNE are all effective dimensionality reduction techniques for micro-expression recognition. The choice of technique depends on the specific requirements of the application, such as computational efficiency or visualization needs. Further research can be conducted to explore the combination of these techniques or the use of other dimensionality reduction methods to improve micro-expression recognition performance.

## VII. FUTURE WORK

The current work includes several areas for future research and improvement of the current work. The key points are:

- Explore combining dimensionality reduction techniques such as RBM, KPCA, and t-SNE to determine whether a hybrid approach can improve micro-expression recognition performance over individual techniques.

- Test other dimensionality reduction techniques, such as autoencoders and UMAP, to assess whether they yield better results than the techniques studied.

- Deep learning models such as Convolutional Neural Networks (CNNs) have shown good performance in facial expression recognition tasks, so Incorporate them after dimensionality reduction.

- Explore micro-expression recognition applications such as lie detection, psychological analysis, and diagnosis of mental illness, where subtle facial expressions are important.

Future extensions will improve the comprehension of micro-expression recognition and its practical applications.

## REFERENCES

[1] Ekman, P. (2003) Darwin, deception, and facial expression, Annals of the New York Academy of Sciences,1000(1), pp. 205–221.

[2] Matsumoto, D., Hwang, H. (2011) Evidence for training the ability to read microexpressions of emotion, Motivation and Emotion, pp. 1–11.

[3] Information for the Dlib library be referred to: http://dlib.net/ .

[4] Information for the CASMEII database will be found on the link: http://casme.psych.ac.cn/casme/e2 .

[5] Lui, R. and Gillies, D. (2016), Overfitting in linear feature extraction for classification of highdimensional image data, Patter Recognition, 53, pp. 73-86.

[6] Bossaerts, P., Carsten, M., Computational Complexity and Human Decision-Making, Trends in Cognitive Sciences, 2017, Vol.21, pp.917-929, doi.org/10.1016/j.tics.2017.09.005 .

[7] Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. Science, 313(5786), 504-507.

[8] G. Hinton, Y. Bengio, Visualizing data using t-SNE, Cost-sensitive Machine Learning for Information Retrieval 33, 2008.

[9] Adegun, I., Vadapalli., H. (2016), Automatic Recognition of Micro-expressions using Local Binary Patterns on Three Orthogonal Planes and

Extreme Learning Machine, 2016 Pattern Recognition Association of South Africa and Robotics and Mechatronics, International Conference (PRASA-RobMech) Stellenbosch, South Africa.

[10] Wang, Y., & Ji, Q. (2016). A Comparative Study of RBM, KPCA and t-SNE for Micro-expression Recognition. In 2016 IEEE 8th International Conference on Biometrics Theory, Applications and Systems (BTAS) (pp. 1-6). IEEE.

[11] Ying, W., Lianghua, H., Pengfei, S. (2012), Face recognition using difference vector plus KPCA, Digital Signal Processing, 22(1), pp. 140-146.

[12] Yan, W., Li, L., Wang, S., & Chen, L. (2013). CASME: A database for spontaneous macro-expression and micro-expression spotting and recognition. In 2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG) (pp. 1-5). IEEE.

[13] Wang, S., Yan, W., Sun, T., Zhao, G., Fu, X. (2016), Sparse tensor canonical correlation analysis for micro-expression recognition, Neurocomputing, 214, pp. 218-232.

[14] García. A., Huerta. A., Zlotnik. S., Díez. P., (2020), A kernel Principal Component Analysis (KPCA) Digest with a New Backward Mapping (pre-image reconstruction), doi.org/10.21203/rs.3.rs-126052/v1.

[15] Maaten, L. V. D., & Hinton, G. (2008). Visualizing data using t-SNE. Journal of Machine Learning Research, 9(Nov), 2579-2605.

[16] Wang, Q., 2012. Kernel principal component analysis and its applications in face recognition and active shape models. arXiv preprint arXiv:1207.3538

[17] W.-J. Yan, X. Li, S.-J. Wang, G. Zhao, Y.-J. Liu, Y.-H. Chen, X. Fu, CASME II: An improved spontaneous micro-expression database and the baseline evaluation, PloS one 9 (1) (2014) e86041.

[18] Bu, Y., Zhao, G., Luo, A., Pan, J., Chen, Y., Restricted Boltzmann Machine: A Non-Linear Substitute for PCA in Spectral Processing, Astronomy&Astrophysics,2015,doi.org/10.1051/0004-6361/201424194.

[19] Cohn, J.F. (2006), Foundations of Human Computing: Facial Expression and Emotion, Int. Conf. on Multimodal Interfaces, 233-238.

[20] Van der Maaten, L.; Hinton, G. Visualizing data using t-SNE. Journal of Machine Learning Research. 2008, 9, 2579–2605

[21] Zhao, G., & Pietikainen, M. (2007). Dynamic texture recognition using local binary patterns with an application to facial expressions. IEEE Transactions on Pattern Analysis and Machine Intelligence, 29(6), 915-928.

[22] Ruiz-Hernandez JA, Pietikäinen M (2013) Encoding Local Binary Patterns Using the Re-Parametrization of the Second Order Gaussian Jet. 10th Proc Int Conf Autom Face Gesture Recognit (FG2013). Shanghai, China: IEEE. doi.org/10.1109/FG.2013.6553709.

[23] Hossin, M.; Sulaiman, M. A review on evaluation metrics for data classification evaluations. Int. J. Data Min. Knowl. Manage. Process 2015, 5, 1–11.

[24] Li, X., Hong,X ., Moilanen, A., Huang, X., Pfister, T., Zhao, G., Pietikainen, M. (2017), Towards Reading Hidden Emotions: A Comparative Study of Spontaneous 3 Micro- Expression Spotting and Recognition Methods, IEEE TRANSACTIONS ON AFFECTIVE COMPUTING, 8, pp.1-15.

[25] Muna, N., Rosiani, U, Yuniarno, E., Purnomo, M. (2017), Subpixel Subtle Motion Estimation of Micro-Expressions Multiclass Classification, 2017 IEEE 2nd International Conference on Signal and Image Processing.

[26] Owayjan, M., Kashour, A., Al Haddad, N., Fadel, M., Al Souki, G. (2012), The Design and Development of a Lie Detection System using Facial Micro-Expressions, 2012 2nd International Conference on Advances in Computational Tools for Engineering Applications (ACTEA), pp. 33 - 38.

[27] Döring, Matthias. Dimensionality Reduction in Machine Learning,

[28] Data Science Blog, 2018, www.datascienceblog.net/post/machine-learning/dimensionality-reduction/.

# A Method for Extracting Traffic Parameters from Drone Videos to Assist Car-Following Modeling

Xiangzhou Zhang, Zhongke Shi

School of Automation, Northwestern Polytechnical University, Xi'an, China

*Abstract*—A new method for extracting traffic parameters from UAV videos to assist in establishing a car-following model is proposed in this paper. The improved ShuffleNet network and GSConv module were introduced into the Yolov7-tiny neural network model as the target detection stage. HOG features and IOU motion metrics are introduced into the DeepSort multi-object tracking algorithm as the tracking matching stage. By building a self-built UAV aerial traffic data set, experiments were conducted to prove that the new method improved a few detection and tracking indicators. In addition, it improves the false detection, missed detection, wrong ID conversion and other phenomena of the previous algorithm, and improves the accuracy and lightweight of multi-target tracking. Finally, gray correlation was applied to analyze the traffic parameters extracted by the new method, and the driver's visual perception of collision was introduced into the car-following model. Through stability analysis, small disturbance simulation and collision risk assessment, the newly proposed traffic flow parameter extraction method has been proven to improve the dynamic characteristics and safety of the car-following model, and can be used to alleviate traffic congestion and improve driving safety.

*Keywords—UAV; Yolov7-tiny; DeepSor; Car-following model; Stability analysis; Traffic congestion; safety assessment*

## I. INTRODUCTION

Traffic congestion leads to low transportation efficiency and air pollution. Frequent traffic accidents lead to casualties and economic losses. These problems are all important challenges facing the development of modern transportation. Researchers have developed numerous models based on traffic flow theory to explain traffic phenomena, thereby improving traffic efficiency and driving safety. Furthermore, the collection of data required for modeling has always been the basis and hot spot of research. UAV aerial video data is extremely informative both in content and time. The UAV is rapidly popularized due to its lightweight, easy operation, and low cost, making them increasingly important in the field of target detection and tracking [1-3]. It is often used in traffic law enforcement and monitoring in various countries. However, how to use drone aerial photography to assist in building a driver behavior model still needs to be explored.

The key to the application of UAV aerial photography data collection lies in the video vehicle detection and tracking algorithm, which extracts its speed, trajectory and other information through vehicle position information at different times. With the rapid development of deep learning technology in recent years, in terms of vehicle target detection, researchers have proposed a variety of improved target detection neural networks for different scenarios and tasks. Among them, the

yolo series of single-stage multi-target detection algorithms is widely used due to its obvious advantages. Makarov et al. [4] used the yolo V2 network to realize the recognition of cars, large vehicles and other objects from the UAV perspective. Hoslain [5] and others migrated YOLO V3 and SSD to the edge-side onboard GPU Jetson TX2. Jetson Xavier implemented UAV detection of vehicles and provided accurate target locations and vehicle types. In addition, the problems introduced by the drone aerial photography perspective have been optimized. For example, to address the problem of an increase in small targets caused by UAVs. Zhang et al. [6] inserted three Spatial Pyramid Pooling modules between the fifth and sixth convolutional layers in front of the three detection heads of the YOLO V3 network to design the silm-yolo V3-SPP3 network. In order to solve the problem of low detection efficiency caused by the sparse and uneven distribution of target categories from the perspective of a drone. Li et al. [7] proposed DS YOLO V3, which added multiple detection heads connected to different layers of the backbone network to detect targets of different sizes. In addition, a multi-scale channel attention fusion module is designed to utilize complementary channel information.

In terms of UAV aerial photography target tracking, commonly used methods based on target trajectory mainly include Karman filtering and deepsort framework. Luo et al. [8] used yolov5 for feature extraction, Kalman filter to extract target motion information and update predictions, and Hungarian matching algorithm to obtain tracking results. Khalkhali et al. [9] proposed SAIKF (Situation Assessment Interactive Kalman Filter), which uses situation assessment information extracted from the traffic history of the same environment to improve tracking performance. The target trajectory prediction based on deepsort is as follows. Ning et al. [10] used yolov5 to obtain the real-time position of the target, and combined with the deepsort framework to achieve the speed measurement of the target. In addition to the above applications, many scholars have made various corresponding improvements to address the problems that arise in multi-target tracking from the UAV perspective. Du et al. [11] used OSNet to replace the simple feature extractor in Deep-SORT, used global clues to associate it with the trajectory, and proposed the EMA (Exponential Moving Average) strategy to achieve a more accurate association between small trajectories and detection results. Huang et al. [12] generated target bounding boxes through different prediction networks, performed cascade matching on all trajectories and detection results, performed unmatched tracking and detection through GIOU matching, and generated the final trajectory.

By establishing a mathematical model of the relative motion relationship between vehicles, car-following models are often suitable for the development of autonomous driving systems and traffic flow simulations. Therefore, it has been a hotspot in traffic flow research, and scholars have achieved rich achievements. Zhang [13] proposed a bi-directional visual angle car-following model considering collision sensitivity, which improved the dynamic characteristics and driving safety of car-following and lane-changing in the traffic flow. Zhang [14] proposed a small-radius curve following model that considers the driver's desired visual angle based on the impact of two-point preview steering decisions and parking sight distance on small-radius curve following behavior from the perspective of the driver's visual characteristics. Liu et al. [15] optimized traffic at signalized intersections by incorporating short-term driving memory on driver behavior. Ma et al. [16] introduced memory effect of headway changes into the driver behavior. Simulations demonstrate it has a significant effect on alleviating traffic congestion.

To sum up, the target detection and tracking technology of drone aerial videos is relatively mature. However, there are still some limitations in using UAV aerial photography to collect information to assist in driving behavior modeling: (1) Due to the size limitations of UAV equipment, there is still room for improvement in high-speed and high-precision traffic information collection. (2) Vehicle data that can quantify the driver's physiological and psychological behavior from a microscopic perspective is difficult to extract. (3) It is difficult to use traffic survey data to model driving behavior while optimizing the dynamic characteristics and safety of traffic flow. This study proposes a new method of assisted modeling of UAV aerial car-following images to solve the above problems. The method we proposed for the above problems has the following advantages: (1) The improved lightweight video detection network and improved target tracking method are conducive to improving the accuracy and speed of traffic information collection in UAV aerial videos. (2) A method that can collect the driver's psychological following behavior is proposed. (3) A method is proposed to use collected traffic data to model driving behavior to simultaneously improve dynamic characteristics and safety.

The remaining parts of this study are structured as follows: Our proposed method for extracting traffic parameters from UAV videos is described in detail in Section II. The improvements of the new traffic parameter extraction method in target detection, tracking and car-following behavior modeling are verified through experiments. The experiment results are described in the Section III. The collected traffic data is used to model the driver's psychological car-following behavior, and its dynamic characteristics and safety are improved in the Section IV. The work of this research is summarized and future work is prospected in Section V.

## II. MATTER AND METHODS

Fig. 1 shows a new method of extracting traffic parameters from UAV video to assist in establishing a car-following model. The method is divided into four stages: image processing, target detection, target tracking and traffic flow modeling. We introduce them separately below.

### A. Transformation of Coordinate Systems

The real traffic parameters are obtained by converting the image coordinates into real coordinates to calculate the traffic parameters. Establishing a transformation matrix between image coordinates and world coordinates is the basis for traffic parameter extraction. Therefore, one frame of image in each video is selected as the reference frame $N_0$ of the video; several marker points are marked on the reference frame, their image coordinates are recorded, and the world coordinates of the same marker points are obtained. The same two sequences in the two coordinate systems have a corresponding relationship. The conversion between the two sets of coordinate systems can be achieved through the perspective projection matrix $T$ of the reference frame image coordinate system and the world coordinate system. The coordinate relationship is as follows:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \boldsymbol{T} \times \begin{bmatrix} x_0 \\ y_0 \\ 1 \end{bmatrix} = \begin{bmatrix} A & D & G \\ B & E & H \\ 0 & 0 & I \end{bmatrix} \times \begin{bmatrix} x_0 \\ y_0 \\ z_0 \end{bmatrix} \qquad (1)$$



Fig. 1. UAV aerial images assist in establishing car-following models system.

where, $[u \ v \ 1]$ is the world coordinate, and $[x_0 \ y_0 \ 1]$ is the image coordinate of the reference frame. Since the image coordinates are two-dimensional coordinates, the coordinate system conversion problem is transformed into a single plane, so the elevation direction parameter in the three-dimensional coordinates is 1. Substituting the world coordinates and the image coordinates of the reference frame into the Eq. (1), the perspective projection matrix $T$ can be obtained.

### B. Image Matching

In actual scenes, the camera may be subject to external interference (such as breeze) and undergo slight displacements and changes in pitch angle. Especially for drones, it is impossible to be completely still during the process of collecting data, and the image coordinates of the same fixed point in different frames will change. Therefore, after completing the conversion of image coordinates and world coordinates, it is necessary to obtain the rotation displacement matrix $w$ in the following formula to calibrate the correspondence between the $n$ th frame and the reference frame $w$ .

$$\begin{bmatrix} x_n \\ y_n \\ 1 \end{bmatrix} = W \times \begin{bmatrix} x_0 \\ y_0 \\ 1 \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} x_0 \\ y_0 \\ 1 \end{bmatrix} \qquad (2)$$

where, $(x_n, y_n)$ is the image coordinate of the nth frame; $(x_0, y_0)$ is the image coordinate of the reference frame $N_0$ ,

### C. Improved Target Detection Algorithm

In the target detection stage, the improved yolov7-tiny target detection model is first used to identify the input image of each frame of the video, and the detecting bounding box is obtained by screening. Then, the detecting bounding box is extracted through the convolutional neural network model to obtain depth features and manual HOG features. These two features are called the appearance features of the detecting bounding box. Finally, the appearance characteristics of the detecting bounding box and the position information characterizing its motion characteristics are input into the tracking and matching stage. This article considers using a more lightweight feature aggregation scheme, which requires fewer parameters and less calculations while ensuring rich features. Therefore, an improved yolov7-tiny target detection model is proposed. First, we use the idea of ShuffleNet, a lightweight network for image classification, to improve Backbone to reduce dense connections and increase network depth. Reducing dense connections can reduce the overall calculation amount, while appropriately increasing the network depth can obtain richer features. Secondly, in the Neck part of the model, the lightweight module GSConv is used for feature aggregation and ELAN is improved to further reduce the amount of model parameters, calculation amount and size while ensuring that rich features are not lost. In this way, the improved YOLOv7-tiny network is shown in Fig. 2.

*1) Improved shuffleNet network:* In the two basic modules of the ShuffleNet [17] network, after the input features enter the right branch, Grouped Convolution (GConv) is first performed, then channel shuffle is performed, and then Depthwise Convolution (DW Conv) with a convolution kernel size of 3 is performed. Finally, perform a Grouped Convolution. The entire process can significantly reduce the parameters of the network, but the resulting feature map will lose semantic information and also cause some loss in accuracy.

Based on the above analysis, the right branches of the two basic modules in Fig. 3(a) are improved. First, depthwise separable convolutions are replaced by group convolution modules. Through grouping, on the premise of increasing the amount of parameters, a certain amount of information can be exchanged between each channel of the feature map; secondly, the channel shuffling operation is changed to a standard convolution with a convolution kernel size of 1, and places it at the end of the branch. Standard convolution operations can play the same role, while also further enriching the semantic information of the feature map without adding parameters.



Fig. 2. Improved yolov7-tiny.

Fig. 3.   Improved shufflenet.

*2) GSConv module:* The GSConv[18] structure is shown in Fig. 4. The input and output channel numbers are C1 and C2 respectively. First, the number of channels is reduced to C2/2 due to a standard convolution, and then it goes through a depth-separable convolution, and the number of channels remains unchanged. The channel information is evenly disrupted by shuffling, and the extracted semantic information is enhanced. Then the expressive ability of image features is improved and enhanced with the fusion of feature information.

When the network performs feature fusion at the Neck layer, semantic information will also be continuously transmitted downward. When the height, width and number of channels of the feature map are continuously compressed and expanded, the loss of part of the semantic information will affect the final prediction. This article introduces the GSConv module into the Neck layer of the network, using the GSConv module instead of standard convolution for upsampling and downsampling, reducing the amount of parameters and calculations of the model, and ensuring the sampling effect to the greatest extent. In addition, GSConv is also introduced into the ELAN module for improvement as shown in Fig. 4. The two convolutions before the Concat layer use the GSConv module to reduce the number of parameters of the model while ensuring detection accuracy and "slimming down" the ELAN module.



Fig. 4.   GSConv module.

### D.  Improved Target Tracking Algorithm

The tracking matching stage obtains motion features and appearance features based on the target detection stage, and measures the similarity with the predicted motion features and appearance features of the existing tracked target. Then the correlation matrix is fused and cascade matching is performed.

The detecting bounding box and tracking target that failed to complete the matching are matched again by IOU. Through two matching processes, the pairing of the current frame detecting bounding box and the existing tracking target is completed. Finally, the status flag of the tracking target in the tracker is updated, and the update and prediction of the Kalman filter are completed.

*1) Motion feature measurement based on IOU:* In order to better measure the motion characteristics in multi-target tracking, this paper proposes to use the IOU measurement of the detecting bounding box and the tracking target prediction frame to replace the Mahalanobis distance measurement method in the DeepSORT algorithm.

In DeepSort, the degree of matching between the detecting bounding box and the motion characteristics of the previously tracked target's predicted position is characterized by the Mahalanobis distance between the two. The Mahalanobis distance $d'(i,j)$ between the tracking target $i$ and the detection target $j$ is expressed as,

$$d'(i,j) = \left(x_j - y_i\right)^{\mathrm{T}} S_i^{-1} \left(x_j - y_i\right) \qquad (3)$$

where, $x_j$ and $y_i$ respectively are the observed quantities for detecting target $j$ and the predicted quantities for tracking the motion of target $i$, and $S_i$ represents the state covariance matrix of the Kalman filter. The larger the Mahalanobis distance, the greater the difference in motion characteristics between the two. False correlations can be excluded by setting a threshold for the Mahalanobis distance.

However, this Mahalanobis distance measure only uses the distance relationship between the detecting bounding box and the tracking target, and cannot accurately describe the motion information of the two. When two targets with similar appearance characteristics are close to each other, it is easy to cause ID switching problems. At the same time, when the tracking target is After occlusion for a period, the Mahalanobis distance is affected by the increased uncertainty of Kalman filter prediction, making it difficult to reliably measure the motion state. The motion feature measurement based on IOU can more accurately describe the positional relationship between the detecting bounding box and the tracking target. The formula is:

$$IOU(i,j) = \frac{Area\left(P_i \cap Q_j\right)}{Area\left(P_i \cup Q_j\right)} \qquad (4)$$

The value range of $IOU(i,j)$ is [0,1], and the IOU distance between the tracking target $i$ and the detection target $j$ is

$$d_{\mathrm{IOU}}(i,j) = 1 - IOU(i,j) \qquad (5)$$

Through the IOU measurement method, the difference in motion features between the detecting bounding box and the tracking target can be more accurately measured, and the

motion feature value threshold matrix expression can be obtained:

$$b_{\text{IOU}}(i,j) = I\left[d_{\text{IOU}}(i,j), t_{\text{IOU}}\right] \quad (6)$$

In the formula, *I* is the indicator function, which takes 1 when the conditions are met, and 0 otherwise. After experiments, $t_{\text{IOU}}$ is set to 0.9, which can eliminate most of the erroneous correlations between the detecting bounding box and the tracking target prediction frame, improve tracking accuracy, and reduce the number of ID switching times.

*2) Appearance feature measurement based on HOG feature:* In order to more accurately associate the matching detecting bounding box and the tracking target in the tracking and matching stage, this paper proposes to fuse the HOG feature distance in the cascade matching view feature measurement. The HOG feature is characterized by statistics and calculation of the gradient direction histogram of the local area of the image. It is often used to characterize the edge information of objects and is widely used in image recognition. It can maintain good invariance to geometric deformation and illumination changes of the detection frame. Moreover, extracting the HOG features of the calculated image requires a small amount of calculation, fast operation, and has little impact on speed performance. By introducing HOG feature distance fusion, it can better reflect the shallow image features of the target and improve the accuracy and robustness of appearance similarity measurement.

In the target detection stage, the size of the obtained detection frame is adjusted to 256x128, the pre-trained ResNet-18 network is input, and the 512-dimensional depth feature $D_j$ is output. Then the HOG features of the detection frame are extracted and the 8505-dimensional HOG feature $H_j$ is output.

The HOG feature distance uses the minimum cosine distance as the measurement criterion, and are only calculated for the features in the $N$ frames closest to the tracking target $i$. They can be expressed by the following formulas respectively.

$$d_{\text{HOG}}(i,j) = \min\left\{1 - H_j^{\text{T}} H_k^i \mid H_k^i \in R_{\text{HOG}}^i\right\}, R_{\text{HOG}}^i = \left\{H_k^i\right\}_{k=1}^N \quad (7)$$

$$R_{\text{HOG}}^i = \left\{H_k^i\right\}_{k=1}^N \quad (8)$$

Similarly, the appearance feature distance measurement also has a threshold $t_F$ to avoid false matching. Its expression is:

$$b_F(i,j) = I\left[d_F(i,j), t_F\right] \quad (9)$$

where, $t_F$ is usually set to 0.2. Combining the two thresholds of Eq. (5) and Eq. (7), the threshold function $b(i,j)$ is obtained, and its formula is as follows:

$$b(i,j) = b_{\text{IOU}}(i,j) \circ b_F(i,j) \quad (10)$$

In the formula, ∘ represents the Hadamard product of the matrix. The fusion coupling matrix $C_{i,j}$ can be obtained from the fusion appearance distance, and the expression is:

$$C_{i,j} = d_F(i,j) \circ b(i,j) \quad (11)$$

The fusion correlation matrix is matched using the Hungarian matching rule to obtain the correct correspondence between the tracking target and the detecting bounding box.

*E. Gray Correlation Calculation*

We use gray correlation to analyze the correlation between the driver's visual factors and acceleration decisions during car-following. Gray relational analysis is well suited for studies of small sample size data sets and fleeting microscopic driving behaviors. The calculation steps are as follows.

*1) Determine the analysis sequence:* We select the acceleration of the following vehicle as the reference sequence, $A = \{A(j) \mid j = 1, 2 \cdots, n\}$; The newly introduced visual angle related parameters are selected as the comparison sequence, $\theta_i = \{\theta_i(j) \mid j = 1, 2, \cdots, n\}$.

*2) Dimensionless variables*

$$x_i(j) = X_i(j) \Big/ \sqrt{\sum_{i=1}^m X_i^2(j)}, j = 1, 2, \ldots, n; i = 0, 1, 2, \ldots, m \quad (12)$$

*3) Calculate gray correlation coefficient*

The gray correlation coefficient between $a_i(j)$ and $\theta_i(j)$ is

$$\xi_i(j) = \frac{\min_i \min_j |a(j) - \theta_i(j)| + \rho \max_i \max_j |a(j) - \theta_i(j)|}{|a(j) - \theta_i(j)| + \rho \max_i \max_j |a(j) - \theta_i(j)|} \quad (13)$$

Let $\Delta_i(j) = |y(j) - x_i(j)|$, the following can be obtained.

$$\xi_i(j) = \frac{\min_i \min_j \Delta_i(j) + \rho \max_i \max_j \Delta_i(j)}{\Delta_i(j) + \rho \max_i \max_j \Delta_i(j)} \quad (14)$$

where, $\rho \in (0, \infty)$ is the resolution coefficient. Generally, the value range of $\rho$ is $(0,1)$. Here we set $\rho = 0.5$.

*4) Calculate gray correlation*

$$r_i = \sum_{j=1}^n \xi_i(j) \Big/ n, j = 1, 2, \cdots, n \quad (15)$$

## III. EXPERIMENTAL RESULTS AND DISCUSSIONS

This experiment collected and annotated a video data set containing 12 groups of traffic scenes named UVACAR-MOT. The video data is 30 frames/s, and one frame is extracted every three frames to form a new video sequence. We selected eight groups as training sets and four groups as validation sets. The experimental setup employs the AutoDL cloud computing platform with a Linux system, 32GB of memory, and PyTorch 1.7 for deep learning. The graphics card is NVDIA Quadro

V100 with 32G of video memory. In the target tracking experiment parameter settings, the initial frame is three frames; the maximum threshold distance for IOU matching is 0.7; and the maximum retained frame for lost tracking is 30. The number of frames for calculating appearance characteristics is N=100. In the target detection experiment parameter settings, the detecting bounding box confidence threshold of yolov7-tiny is 0.3, and the IOU threshold of NMS non-maximum filtering in yolov7-tiny is 0.5.

## A. Target Detection

Table I shows that compared with the SSD algorithm and the yolov7-tiny algorithm, the improved yolov7-tiny algorithm has an increase in average accuracy and a decrease in parameter scale. The new algorithm has the characteristics of better detection accuracy, low parameters and low computational load. Compared with the mainstream yolo5s, it is more suitable for deployment on drones. Although the accuracy of the new algorithm is slightly lower than that of yolox-s, the parameters and model size are lower.

TABLE I.      OTHER IMAGE DETECTION ALGORITHMS ON UVACAR-MOT

|  | mAP_0.5% | mAP_0.95/% | Params/$10^6$ | Flops/$10^9$ | Size/MB |
|---|---|---|---|---|---|
| SSD | 49.1 | 29.2 | 40.3 | 371.2 | 267.6 |
| YoloV7-tiny | 54.1 | 36.2 | 5.8 | 14.3 | 11.9 |
| Yolo5s | 55.3 | 35.8 | 6.6 | 15.2 | 13.5 |
| Yolox-s | 58.3 | 39.9 | 8.9 | 25.7 | 67.5 |
| ours | 58.1 | 39.5 | 4.7 | 13.2 | 10.9 |

Fig. 5 and Fig. 6 respectively show the improvement effect of the improved yolov7-tiny on missed detections and false detections. Fig. 5(a) shows the missed detection of car no. 5 under occlusion for the yolov7-tiny algorithm. Fig. 5(b) shows that car no.5 is still detected under the improved yolov7-tiny algorithm. This shows that GSConv in the new algorithm enhances feature richness and makes detection more accurate. Fig. 6(a) shows that yolov7-tiny false detection curbstones as vehicles. The detection accuracy of the improved yolov7-tiny in the same detecting bounding box is shown in Fig. 6(b). This shows that the improved shuffleNet in the new algorithm effectively enhances the feature extraction capability and improves the accuracy of vehicle detection.



(a) yolov7-tiny                    (b) Improved yolov7-tiny

Fig. 5.   Improved target detection algorithm improves missed detections.



(a)yolov7-tiny                    (b) Improved yolov7-tiny

Fig. 6.   Improved target detection algorithm improves false detections.

## B. Target Tracking

Table II shows that most of the indicators of improved yolov7-tiny+deepsort in the UAVCAR-MOT data set are better than the deepsort tracking algorithm. The improvement of specific MOTA and MOTP indicators shows that the tracking accuracy has been greatly improved. IDsw is reduced to two times, which proves that the number of ID switching times for tracking the same target is very small and the tracking retention ability is strong.

TABLE II.      OTHER IMAGE TRACKING ALGORITHMS ON UVACAR-MOT

|  | MOTA/% | MOTP/% | IDsw | MT/% | FPS |
|---|---|---|---|---|---|
| EAMTT | 51.88 | 73.81 | 45 | 71 | 16.77 |
| POI | 64.34 | 70.56 | 43 | 72 | 17.25 |
| Sort | 58.45 | 74.79 | 39 | 70 | 29.92 |
| Deepsort | 61.41 | 74.73 | 34 | 72 | 18.53 |
| Ours | 62.02 | 76.51 | 2 | 74 | 28.83 |



(a) yolov7-tiny+deepsort          (b) yolov7-tiny+deepsort

(c) Improved yolov7-tiny+deepsort    (d) Improved yolov7-tiny+deepsort

Fig. 7.   Improved target tracking algorithm to improve ID hopping.



(a) yolov7-tiny+deepsort          (b) yolov7-tiny+deepsort

(c) Improved yolov7-tiny+deepsort    (d) Improved yolov7-tiny+deepsort

Fig. 8.   Improved target tracking algorithm improves false detection of shadow.

Fig. 7 shows the different tracking effects of the improved front and rear tracking algorithms when blocked by trees. In Fig 7(a), the ID of the same vehicle in the two bounding box s

is switched from 3 to 2 when using the improved yolov7-tiny+deepsort algorithm. However, no such switching occurs after using the yolov7-tiny+deepsort algorithm in Fig. 7(b). This reflects that the improved yolov7-tiny+deepsort algorithm has better tracking capabilities. Fig. 8 shows the impact when light and shadow changes occur. Fig. 8(a) shows the red arrow vehicle being given two tracking IDs and boxes of different sizes. However, such an error did not occur in Fig. 8(b). This shows that the fusion of appearance feature measurement with HOG feature distance can improve the accuracy of appearance feature measurement and tracking accuracy.

*C. Target Tracking*

In the traffic phenomenon investigation experiment, we chose Jinye Road, Xi'an City, China, and its satellite map is shown in Fig. 9. The aerial photography scene is shown in Fig. 10. We chose clear and windless weather at 3 pm. Furthermore, there are more types and numbers of vehicles on this road, but there is no congestion. The image collection equipment uses a zoom drone, and the video quality is 1080p, 30fps. In order to avoid errors caused by image distortion in the later stage as much as possible, we hovered the drone at 90° directly above the road to shoot and kept the height at 145m. Table III shows one of our multiple sets of aerial traffic data. We selected the vehicle with ID 20 as the following car (width=1.6m) and extracted relevant image and driving data.

The correlation between the driver's visual factors and vehicle acceleration in Table IV calculated from Table III are all above 0.65. This shows that they are strongly related. In addition, we found that when the driver makes the acceleration and deceleration decision, the consideration of the possibility of collision is more important than the headway and velocity.



Fig. 9.   Google maps for car-following.



Fig. 10. Target tracking for car-following.

TABLE III.    AERIAL TRAFFIC DATASET

| Pixel coordinates | $a_n(t)$ | $v_n(t)$ | $\theta_n(t)$ | $\dot{\theta}_n(t)$ | $\varphi_m(t)$ | $\dot{\varphi}_m(t)$ | $\dot{\theta}_m(t)/\theta_m(t) - \dot{\beta}_m(t)/\beta_m(t)$ | $\dot{\varphi}_m(t)/\varphi_m(t)$ |
|---|---|---|---|---|---|---|---|---|
| (723,351) | 0.0500 | 12.1040 | -0.0019 | -0.0010 | 0.1294 | 0.0045 | 0.0348 | -0.5498 |
| (724,351) | 0.1750 | 12.2500 | -0.0021 | -0.0012 | 0.1303 | 0.0062 | 0.0479 | -0.5852 |
| (726,352) | -0.0125 | 12.3850 | -0.0024 | -0.0015 | 0.1316 | 0.0084 | 0.0635 | -0.6620 |
| (728,353) | 0.1000 | 12.5082 | -0.0027 | -0.0021 | 0.1332 | 0.0086 | 0.0644 | -0.8065 |
| (730,353) | -0.0875 | 12.4025 | -0.0031 | -0.0028 | 0.1350 | 0.0115 | 0.0854 | -0.9518 |
| (732,354) | -0.0250 | 12.1850 | -0.0036 | -0.0042 | 0.1373 | 0.0089 | 0.0652 | -1.1941 |
| (735,356) | -0.1000 | 11.7400 | -0.0045 | -0.0068 | 0.1391 | 0.0085 | 0.0614 | -1.5910 |
| (737,357) | -0.2000 | 11.3600 | -0.0058 | -0.0088 | 0.1408 | 0.0065 | 0.0459 | -1.5991 |
| (739,356) | -0.1750 | 10.9200 | -0.0076 | -0.0144 | 0.1421 | 0.0048 | 0.0335 | -2.0249 |
| (743,357) | -0.1875 | 10.6650 | -0.0105 | -0.1000 | 0.1430 | 0.0012 | 0.0081 | -10.2236 |

TABLE IV.    GRAY RELATIONAL DEGREE

| factor | $V_n(t)$ | $\theta_n(t)$ | $\dot{\theta}_n(t)$ | $\varphi_m(t)$ | $\dot{\varphi}_m(t)$ | $\dfrac{\dot{\theta}_m(t)}{\theta_m(t)} - \dfrac{\dot{\beta}_m(t)}{\beta_m(t)}$ | $\dfrac{\dot{\varphi}_m(t)}{\varphi_m(t)}$ |
|---|---|---|---|---|---|---|---|
| Correlation | 0.807 | 0.820 | 0.802 | 0.711 | 0.790 | 0.794 | 0.873 |

## IV. APPLICATIONS OF MODELING

### A. Baseline Model

Based on the optimized speed following model, Jiang [19] proposed a classic full speed difference model with over a thousand references.

$$a_n(t) = \alpha \left\{ V\left[\Delta x_n(t)\right] - v_n(t) \right\} - \lambda \Delta v_n(t) \qquad (16)$$

where, $\alpha$ and $\lambda$ are the sensitivity coefficient. $\Delta x_n(t)$ and $\Delta v_n(t)$ are respectively the relative distance and velocity between the front and rear vehicles.

However, from the perspective of driver psychology, the most important perceptual information in car following behavior may be visual information. so Jin [20] replaces the traditional headway with the visual angle in following behavior.

$$a_n(t) = \alpha \left\{ V\left[\theta_n(t)\right] - v_n(t) \right\} - \lambda \, d\theta_n(t)/dt \qquad (17)$$

$$\theta_n(t) = w/\left(\Delta x_n(t) - l\right) \qquad (18)$$

where, $\theta_n(t)$ is the visual angle of drivers. $w$ and $l$ are the width and length of leading vehicles. $d\theta_n(t)/dt$ represents change rate of visual angles. $V(\theta_n)$ is the optimized velocity.

### B. New Model Derivation

The TTC is the time until vehicles crash assuming the collision path and velocity differential are maintained [21]. The driver of the following car will take control measures such as acceleration or deceleration according to the change in TTC with the leading vehicle [22]. The commonly used traffic accident alternative evaluation index TTC is used as our reference. We combine the correlation between the driver's visual psychology and driving operations in Table 4 to try to establish a more realistic and accurate car-following model to quantify the risk perception of collision accidents. Therefore, driver sensitivity to lateral collision is incorporated into the visual angle model to improve the stability and safety of traffic flow. The TTC expression for a single lane can be written as

$$\text{TTC}_m(t) = \frac{x_{m+1}(t) - x_m(t)}{v_m(t) - v_{m+1}(t)} = \frac{\varphi_m}{\dot{\varphi}_m} \qquad (19)$$

where, $x_{m+1}(t)$, $x_m(t)$ and $v_{m+1}(t)$, $v_{m+1}(t)$ are the position and velocity of vehicles. However, the car-following considering lateral influence has the following geometric relationship as shown in Fig. 11.



Fig. 11. Car-following behavior with visual angle.

$$S_m(t) = \sin \beta_m(t) \cdot \frac{L_m(t)}{\sin \gamma_m(t)} \qquad (20)$$

$$\dot{S}_m(t) = \frac{\left(\dot{L}_m(t)\sin\beta_m(t) + L_m(t)\cos\beta_m(t)\dot{\beta}_m(t)\right)}{\sin\gamma_m(t)} \qquad (21)$$

where, $L_m(t)$ is the distance between the midpoint of tail of preceding vehicles and the midpoint of head of the following vehicles, $S_m(t)$ is the distance between the head of preceding vehicles and the conflict point. The angle $\beta_m(t)$ defined by $L_m(t)$ and the distance to the collision point. We assume that the leading car maintains constant velocity and $\varphi_m(t)$ for the current brief time. we roughly think that $\tan\beta_m(t) = \beta_m(t)$. Substituting Eq. (20) into Eq.(21) and eliminating the $\sin\gamma_m(t)$.

$$\frac{\dot{S}_m(t)}{S_m(t)} = \frac{\dot{L}_m(t)}{L_m(t)} + \frac{\dot{\beta}_m(t)}{\beta_m(t)} \qquad (22)$$

Then the potential collision with side vehicle can be expressed.

$$\frac{1}{\text{TTC}} = \frac{\dot{\theta}_m(t)}{\theta_m(t)} - \frac{\dot{\beta}_m(t)}{\beta_m(t)} \qquad (23)$$

Therefore, we extend the TTC indicator in car-following behavior to car-following behavior, and can obtain a new visual angle model considering collision visual sensitivity:

$$a_m(t) = \alpha \left\{ V\left[\theta_m(t), \varphi_m(t)\right] - v_m(t) \right\} - \lambda \left[ (1 - \xi_m)\dot{\theta}_m(t) \right.$$
$$\left. + \xi_m\dot{\varphi}_m(t) \right] - \kappa \left[ (1 - \xi_m)\left(\frac{\dot{\theta}_m(t)}{\theta_m(t)} - \frac{\dot{\beta}_m(t)}{\beta_m(t)}\right) + \xi_m\frac{\dot{\varphi}_m(t)}{\varphi_m(t)} \right] \qquad (24)$$

The comprehensive optimization velocity is as follows:

$$V\left[\theta_m(t), \varphi_m(t)\right] = V_1 + V_2 \tanh\left\{ C_1\left[(1 - \xi_m)w/2\tan\left(\theta_m(t)/2\right)\right.\right.$$
$$\left.\left. + \xi_m w/2\tan\left(\varphi_m(t)/2\right)\right] - C_2 \right\} \qquad (24)$$

These parameters $V_1, V_2, C_1, C_2$ were verified by Zhang [14].

### C. Stability Analysis

We assume that same size vehicles traveling on a ring road with a uniform flow as the initial state. Each vehicle maintains the same headway $h$ with adjacent vehicles at a uniform velocity $V(h)$. So, the initial moment can be considered as,

$$x_m^0(t) = hn + V(\theta_0, \varphi_0)t \qquad (25)$$

where $\theta_0 = 2\arctan\left[w/2(h - l)\right]$, $\varphi_0 = 2\arctan\left[w/(4h - 2l)\right]$.

When the perturbation $y_m(t)$ appear, we can obtain

$$x_m(t) = x_m^0(t) + y_m(t) \qquad (26)$$

Visual angles formed by the driver when he observes himself and the side vehicle, and angles formed by side vehicle and the collision position are respectively expressed as,

$$\theta_m(t) = 2\arctan\left[ w / 2\left(\Delta x_{m,m+1}(t) - l\right) \right] \qquad (27)$$

$$\varphi_m(t) = 2\arctan\left[ w / 2\left(\Delta x_{m,m+2}(t) - l\right) \right] \qquad (28)$$

$$\beta_m(t) = \arctan\left[ d\tan\gamma / \left(\Delta x_{m,m+1}(t) - l\right) - \tan\gamma \right] \qquad (29)$$

To facilitate calculation of the visual angle expression is linearized using Taylor expansion Substituting Eq. (27) into Eq. (27), Eq. (28) and Eq. (29) higher-order terms of $\Delta y_m(t)$ can be rounded off. Drivers' visual angles are expressed as follows

$$\theta_m(t) = 2\arctan\frac{w}{2(h-l)} - \frac{w}{(h-l)^2 + (w/2)^2}\Delta y_{m,m+1}(t) \quad (30)$$

$$\beta_m(t) = \arctan\left(\frac{d\tan\gamma}{h-l} - \tan\gamma\right) \\ - \frac{d\tan\gamma}{(h-l)^2 + \tan^2\gamma(d-h+l)^2}\Delta y_{m,m+1}(t) \qquad (31)$$

$$\varphi_m(t) = 2\arctan\frac{w}{2(h-l)} - \frac{w}{(h-l)^2 + (w/2)^2}\Delta y_{m,m+2}(t) \quad (32)$$

Substituting (31)~ (33) into (24), the following is obtained

$$\ddot{y}_m(t) = \alpha A\dot{V}(\theta_0,\varphi_0)\left\{(1-\xi_m)\Delta y_{m,m+1}(t) - \xi_m\Delta y_{m,m+2}(t)\right. \\ \left. - \dot{y}(t)\right\} + \lambda A\left((1-\xi_m)\Delta\dot{y}_{m,m+1}(t) + \xi_m\Delta\dot{y}_{m,m+2}(t)\right) \qquad (33)$$

$$+\kappa\left((1-\xi_m)\left(\frac{-\Delta\dot{y}_{m,m+1}(t)}{B - \Delta\dot{y}_{m,m+1}(t)} + \frac{\Delta\dot{y}_{m,m+1}(t)}{C - \Delta\dot{y}_{m,m+1}(t)}\right)\right.$$

$$\left. + \xi_m\frac{-\Delta\dot{y}_{m,m+2}(t)}{B - \Delta\dot{y}_{m,m+2}(t)}\right)$$

where,

$$\begin{cases} A = -\dfrac{w}{(h-l)^2 + (w/2)^2} \\[2mm] B = \dfrac{(h-l)^2 + (w/2)^2}{w}2\arctan\dfrac{w}{2(h-l)} \\[2mm] C = \dfrac{(h-l)^2 + \tan^2\gamma(d-h+l)^2}{d\tan\gamma}\arctan\left(\dfrac{d\tan\gamma}{h-l} - \tan\gamma\right) \end{cases} \qquad (34)$$

where, $\dot{V}(\theta_0,\varphi_0) = dV(\theta,\varphi)/d(\theta)\big|_{\theta=\theta_0,\varphi=\varphi_0}$ ,We expand $\Delta y_m(t) = B\exp(ikm+zt)$ by Fourier series and substitute it into Eq. (34).

$$z^2 = \alpha A\dot{V}(\theta_0,\varphi_0)\left\{(1-\xi_m)\left(e^{ik}-1\right) - \xi_m\left(e^{2ik}-1\right)\right. \\ \left. - \dot{y}_n(t)\right\} + \lambda A\left((1-\xi_{m+1})z\left(e^{ik}-1\right) + \xi_m z\left(e^{2ik}-1\right)\right) \\ - \kappa\left((1-\xi_m)\left(\frac{-z\left(e^{ik}-1\right)}{B-\left(e^{ik}-1\right)} + \frac{z\left(e^{ik}-1\right)}{C-\left(e^{ik}-1\right)}\right)\right. \\ \left. + \xi_m\frac{-z\left(e^{2ik}-1\right)}{B-\left(e^{2ik}-1\right)}\right) \qquad (35)$$

Expanding $z$ according to $z = z_1(ik) + z_2(ik)^2 + \cdots$, we can obtain as follows:

$$\begin{cases} z_1 = (1-\xi_m)A\dot{V} \\[2mm] z_2 = \dfrac{(1+3\xi_m)A\dot{V}}{2} + \dfrac{1}{\alpha}\left[\dfrac{\kappa(1-\xi_m)(B-C) + 2\kappa\xi_m BC}{BC}\right. \\[3mm] \left. + A\lambda(1+\xi_m)\right]z_1 - \dfrac{1}{\alpha}z_1^2 \end{cases} \qquad (36)$$

According to the hypothesis of long-wave expansion, if $z_2$ is positive, traffic flow is still steady under small disturbances. Therefore, the stability curve is obtained.

$$\alpha = \frac{2A(1-\xi_m)^2(\dot{V}+\lambda)}{1+3\xi_m} - \frac{2\kappa(1-\xi_m)^2(B-C) + 4\kappa(1-\xi_m)}{BC} \qquad (39)$$



Fig. 12. Stability curves at different $\kappa$.

The stability curves of the car-following considering visual sensitivity to collision at different $\kappa$ are depicted in Fig. 12. The upper and lower sides of the stability curve respectively are stable and unstable regions. Stability curves of different collision sensitivity coefficients $\kappa$ are depicted in Fig. 12. We set $\lambda=1, \xi=0.1, b=3.6, w=1.6$. When the parameters are $\kappa=0, b=0, \xi=0$, our model degenerates into the VAM and

the following vehicle only optimizes its own velocity according to leading vehicles on the current lane. The area formed by the x-axis and the stability curves of the new model shrinks as $\kappa$ increases, and is smaller than the area of FVDM and VAM in the Fig. 12. The order of areas formed by curves and x-axis in Fig. 12 is FVDM> VAM> SCVAM, and they gradually increase as $\kappa$ decreases. The findings indicate that the platoon is more stable if the visual collision factor is introduced into the car-following model. In addition, the stability gradually improves with the increase of visual collision sensitivity, and the dynamic performance is better than the full velocity difference (FVDM) and visual angle models (VAM).

### D. Simulation

Within this part of the research, we performed a series of simulations on the evolution of small perturbations for the new visual model Eq. (24) to analyze its dynamic performance. To validate the theoretical results obtained above, numerical simulations of Eq. (15) with periodic boundary conditions is given. The time step is $0.1\text{s}$. The initial setting was adopted.

$$\begin{cases} \Delta x_i(0) = \Delta x_i(1) = 5.0, & (i \neq 30,31) \\ \Delta x_i(0) = \Delta x_i(1) = 5.0 - 0.1, & (i = 30) \\ \Delta x_i(0) = \Delta x_i(1) = 5.0 + 0.1, & (i = 31) \end{cases} \quad (37)$$

where 100 vehicles are traveling on a 1500m ring road.



Fig. 13. Space-time graphs of headway at different $\kappa$.



Fig. 14. (a) Hysteresis loops with different $\kappa$ (b) Velocity images of all vehicles with different $\kappa$.

Fig. 13(a) to Fig. 13(d) reflect the temporal and spatial evolution of the car-following model considering the visual sensitivity of collision with different $\kappa$ from 2000s to 2300s. The parameters are $\alpha = 0.41, b = 3.6, \lambda = 20, w = 2$. The vehicle platoon is unstable at $\kappa = 0,10,20,30$ by stability condition Eq.(39). Small disturbances in the system will gradually amplify over time and cause vehicle platoon congestion, which can be clearly seen in Fig. 11(a) to Fig. 11(c). Especially when the lateral distance $b$ and collision sensitivity $\kappa$ are not considered, the newly proposed visual angle car-following model Eq. (24) degenerates to VAM ($\kappa = 0, b = 0$), the greatest fluctuation is in headway as shown in Fig. 11(a). When $\kappa = 30$, the stability condition is satisfied, and the local disturbance of the distance between the vehicles in the platoon will gradually return to stable state in Fig. 11(d). Similar，the areas of hysteresis loops decrease gradually when the visual sensitivity of collision $\kappa$ increase in Fig. 14(a). As the visual sensitivity of collision $\kappa$ increases, the velocity fluctuation decreases are depicted in Fig. 14(b).

Fig. 13 and Fig. 14 show that the stability of the vehicle platoon is well maintained when taking into account the collision visual sensitivity to the lateral and longitudinal leading vehicles. As the sensitivity $\kappa$ increases, the stability gradually increases.

### E. Safety Assessment

The traffic conflict theory can observe a large amount of non-accident data before the accident, and estimate the dangerous action highly related to the accident by analyzing the vehicle. Numerous indicators have been developed for alternative safety measures to evaluate collisions risk. TTC is the most used indicator among them. However, traditional traffic conflict alternative indicators have two shortcomings: 1) They do not start from the driver's real visual perspective. The vehicle is regarded as a point, and the impact of its size on collision is not considered. 2) The potential risk of collision with the two vehicles in front is not considered when following a car. Therefore, we propose a new car-following collision risk indicator STTC, which is expressed as follows.

$$\text{STTC} = (1 - \xi_m)\left(\frac{\dot{\theta}_m(t)}{\theta_m(t)} - \frac{\dot{\beta}_m(t)}{\beta_m(t)}\right) + \xi_m \frac{\dot{\varphi}_m(t)}{\varphi_m(t)} \quad (38)$$

We also use the small perturbation evolution scenario in previous section to study potential collision risks. When the 30th vehicle in the uniform traffic flow suddenly suffered a small disturbance, we selected it and five vehicles at the front and rear, a total of eleven vehicles, for potential collision risk assessment. Since the simulation step in car-following behavior is small and the velocity difference between nearby vehicles is small, we use a larger range of STTC to evaluate the risk. Similar to the common threshold standard of the TTC index, we define $\text{Frequency}_{\text{STTC}<30s}$ as the traffic environment with potential collisions. With reference to the common international standards of TTC, $\text{Frequency}_{\text{STTC}<3s}$ and $\text{Frequency}_{\text{STTC}<5s}$ are defined as the frequency of serious traffic conflicts and relatively serious traffic conflicts.

The safety of the entire platoon is evaluated by the dispersion of velocity and headway. The standard deviation of velocity and headway can be calculated by following formulas.

$$\text{SD}_v = \sum_{t=1}^{T}\sum_{m=1}^{M}\left(v_{t,m}-\bar{v}\right)\Big/\left(M\cdot T-1\right), \quad \text{SD}_h = \sum_{t=1}^{T}\sum_{m=1}^{M}\left(h_{t,m}-\bar{h}\right)\Big/\left(M\cdot T-1\right) \quad (39)$$

where $\bar{v}$ and $\bar{h}$ are the average velocity and headway.



(a) $\kappa = 0, \; b = 0$      (b) $\kappa = 10$

(c) $\kappa = 20$      (d) $\kappa = 30$

Fig. 15. Statistical distributions for STTC with different $\kappa$

TABLE V.     TOTAL OCCURRENCES OF STTC FOR DIFFERENT $\kappa$

| $\kappa$ | $0(\xi=0)$ | 10 | 20 | 30 |
|---|---|---|---|---|
| Frequency$_{\text{STTC}<3s}$ | 720 | 441 | 0 | 0 |
| Frequency$_{\text{STTC}<5s}$ | 952 | 875 | 0 | 0 |
| Frequency$_{\text{STTC}<30s}$ | 2189 | 1646 | 442 | 0 |

TABLE VI.     SD OF HEADWAY AND VELOCITY WITH DIFFERENT $\kappa$

| $\kappa$ | $0\,(\xi=0)$ | 10 | 20 | 30 |
|---|---|---|---|---|
| Velocity | 4.61 | 3.03 | 0.49 | 0.05 |
| Headway | 6.15 | 3.59 | 0.52 | 0.07 |

Fig. 15 shows that when the collision sensitivity is not considered $(\kappa=0, b=0)$, the frequency of $VT<3s$ is the highest, and the risk of serious collision conflict is the highest. In addition, the frequency of dangerous driving within 5s threshold and the potential collision conflict within 30s threshold are also the highest. However, as the driver's sensitivity $\kappa$ to the collision of two vehicles intersecting in front increases, the frequency of the three safety evaluation indicators within the threshold gradually decreases until it completely enters a safe driving state.

The STTC values of different collision sensitivities at different risk thresholds are counted in Table V. The standard deviation of the headway and velocity at different collision sensitivities $\kappa$ are descriptive statistics in Table VI. This shows that increasing the sensitivity to the collision of

horizontal and longitudinal leading vehicles can reduce dispersion of headway and velocity, thereby improving driving safety.

Through the above simulation experiments and statistical analysis, we found that after analyzing the traffic behavior data collected by flexible and convenient drone images, the newly established car-following model can describe the collision risk of driving behavior. And by enhancing the sensitivity of the newly introduced visual collision factor in the model as in Eq.(24), the driver's safety in car following behavior will be improved. This puts forward new ideas for safety modeling and evaluation of driving behavior.

## V. CONCLUSION

To sum up, the method in this article for extracting traffic parameters from UAV video and applying it to assist in establishing a car-following model can improve the accuracy and lightweight of multi-target detection and tracking. The improved ShuffleNet network and GSConv module introduce the Yolov7-tiny target detection stage to reduce the number of parameters and calculations of the model and ensure accuracy. HOG features and IOU motion metrics are introduced into the DeepSort multi-target tracking algorithm to improve the target appearance representation capabilities and the accuracy of tracking targets. The traffic parameters extracted by the new method can be used to analyze driving psychology and car-following behavior, and the analysis results can be used to model car-following behavior with the aim of enhancing both the safety and stability of traffic flow. This will in turn ease traffic congestion and reduce driver collision risks. In the future, we will further improve vehicle tracking in more complex traffic scenarios, and verify the accuracy of traffic parameter extraction through on-board comparison.

### REFERENCES

[1] A. Ammar, A. Koubaa, M. Ahmed, A. Saad and B. Benjdira, "Vehicle Detection from Aerial Images Using Deep Learning: A Comparative Study." Electronics. vol. 10, no. 7. 2021.

[2] H. Yao, R. Qin and X. Chen, "Unmanned Aerial Vehicle for Remote Sensing Applications-A Review." Remote Sensing. vol. 11, no. 12. 2019.

[3] R. Ravindran, M. J. Santora and M. M. Jamali, "Multi-Object Detection and Tracking, Based on DNN, for Autonomous Vehicles: A Review." Ieee Sensors Journal. vol. 21, no. 5, pp. 5668-5677. 2021.

[4] S. B. Makarov, V. A. Pavlov, A. K. Bezborodov, A. I. Bobrovskiy and D. Ge, "Multiple Object Tracking Using Convolutional Neural Network on Aerial Imagery Sequences." 2021.

[5] S. Hossain and D.-J. Lee, "Deep Learning-Based Real-Time Multiple-Object Detection and Tracking from Aerial Imagery via a Flying Robot with GPU-Based Embedded Devices." Sensors. vol. 19, no. 15. 2019.

[6] P. Zhang, Y. Zhong and X. J. I. Li, "SlimYOLOv3: Narrower, Faster and Better for Real-Time UAV Applications." 2019.

[7] Z. Li, X. Liu, Y. Zhao, B. Liu, Z. Huang and R. Hong, "A lightweight multi-scale aggregated model for detecting aerial images captured by

UAVs." Journal of Visual Communication and Image Representation. vol. 77. 2021.

[8] X. Luo, R. Zhao and X. Gao, "Research on UAV Multi-Object Tracking Based on Deep Learning," 2021 IEEE International Conference on Networking, Sensing and Control (ICNSC), Xiamen, China, pp. 1-6, 2021.

[9] M. B. Khalkhali, A. Vahedian and H. S. Yazdi, "Situation Assessment-Augmented Interactive Kalman Filter for Multi-Vehicle Tracking." Ieee Transactions on Intelligent Transportation Systems. vol. 23, no. 4, pp. 3766-3776. 2022.

[10] M. Ning, X. Ma, Y. Lu, S. Calderara and R. Cucchiara, SeeFar: Vehicle Speed Estimation and Flow Analysis from a Moving UAV, 21st International Conference on Image Analysis and Processing (ICIAP), Lecce, ITALY, 2022, pp. 278-289.

[11] Y. Du, J. Wan, Y. Zhao, B. Zhang, Z. Tong, J. Dong and I. C. Soc, GIAOTracker: A comprehensive framework for MCMOT with global information and optimizing strategies in VisDrone 2021, 18th IEEE/CVF International Conference on Computer Vision (ICCV), Electr Network, 2021, pp. 2809-2819.

[12] W. Huang, X. Zhou, M. Dong and H. Xu, "Multiple objects tracking in the UAV system based on hierarchical deep high-resolution network." Multimedia Tools and Applications. vol. 80, no. 9, pp. 13911-13929. 2021.

[13] X. Z. Zhang, Z. K. Shi, J. Z. Chen and L. J. Ma, "A bi-directional visual angle car-following model considering collision sensitivity." Physica a-Statistical Mechanics and Its Applications. vol. 609. 2023.

[14] X. Zhang, Z. Shi, S. Yu and L. Ma, "A new car-following model considering driver's desired visual angle on sharp curves." Physica a-Statistical Mechanics and Its Applications. vol. 615. 2023.

[15] D.-W. Liu, Z.-K. Shi and W.-H. Ai, "Enhanced stability of car-following model upon incorporation of short-term driving memory." Communications in Nonlinear Science and Numerical Simulation. vol. 47, pp. 139-150. 2017.

[16] G. Ma, M. Ma, S. Liang, Y. Wang and Y. Zhang, "An improved car-following model accounting for the time -delayed velocity difference and backward looking effect." Communications in Nonlinear Science and Numerical Simulation. vol. 85. 2020.

[17] W. Sun, B. Fu and Z. Zhang, "Maize Nitrogen Grading Estimation Method Based on UAV Images and an Improved Shufflenet Network." Agronomy-Basel. vol. 13, no. 8. 2023.

[18] X. Zhao and Y. Song, "Improved Ship Detection with YOLOv8 Enhanced with MobileViT and GSConv." Electronics. vol. 12, no. 22. 2023.

[19] R. Jiang, Q. S. Wu and Z. J. Zhu, "Full velocity difference model for a car-following theory." Physical Review E. vol. 64, no. 1. 2001.

[20] S. Jin, D.-H. Wang, Z.-Y. Huang and P.-F. Tao, "Visual angle model for car-following theory." Physica a-Statistical Mechanics and Its Applications. vol. 390, no. 11, pp. 1931-1940. 2011.

[21] M. M. Minderhoud and P. H. Bovy, "Extended time-to-collision measures for road traffic safety assessment." Accident analysis and prevention. vol. 33, no. 1, pp. 89-97. 2001.

[22] D. N. Lee, "A theory of visual control of braking based on information about time-to-collision." Perception. vol. 38, no. 6, pp. A43-A65. 2009.

# A Review of Fake News Detection Techniques for Arabic Language

Taghreed Alotaibi[1], Hmood Al-Dossari[2]

Computer Science Department, Imam Mohammad Ibn Saud Islamic University, Riyadh, Saudi Arabia[1]
Information Systems Department, College of Computer and Information Sciences, King Saud University, Riyadh, Saudi Arabia[1, 2]

*Abstract*—The growing proliferation of social networks provides users worldwide access to vast amounts of information. However, although social media users have benefitted significantly from the rise of various platforms in terms of interacting with others, e.g., expressing their opinions, finding products and services, and checking reviews, it has also raised critical problems, such as the spread of fake news. Spreading fake news not only affects individual citizens but also governments and countries. This situation necessitates the immediate integration of artificial intelligence methodologies to address and alleviate this issue effectively. Researchers in the field have leveraged different techniques to mitigate this problem. However, research in the Arabic language for fake news detection is still in its early stages compared with other languages, such as English. This review paper intends to provide a clear view of Arabic research in the field. In addition, the paper aims to provide other researchers working on solving Arabic fake news detection problems with a better understanding of the common features used in extraction, machine learning, and deep learning algorithms. Moreover, a list of publicly available datasets is provided to give an idea of their characteristics and facilitate researcher access. Furthermore, some of limitations and challenges related to Arabic fake news and rumor detection are discussed to encourage other researchers.

*Keywords*—*Fake news detection; rumors; classification; Arabic language*

## I. INTRODUCTION

The influence of traditional information sources, such as television and newspapers, and how users gather and consume news has diminished in comparison to earlier times. The expansion of social media platforms has been a key factor in this transition. Social media platforms are another kind of technological innovation. Users can access platforms such as Twitter, Facebook, and Instagram to build their profiles, share opinions, interact with others having shared interests, and facilitate cycles of cooperation. Over the years, more users have shown an interest in using social media. According to Kepios, there were 4.76 billion social media users around the world in 2023, which equates to 59.4% of the world's total population [1]. The reason behind this spike in the use of social media is that social media platforms are designed to be more attractive and highly suitable for social communication. Social media has also become a way for companies and governments to reach the people by providing news, showcasing their services, giving updates, or launching marketing campaigns. However, some limitations arise from using social media, such as the spread of fake news and rumors, as well as spamming. Many individuals who use

social media platforms to stay in touch with friends and family also use these platforms to find news and information. According to a report from the Pew Research Center, 48% of adults in the US use social media as a news source [2].

Fake news is a major problem that started early, attracting significant attention in 2016 during the US presidential elections. Different fake news items and rumors usually appear during special events and cover different domains, such as those related to elections, natural disasters, or health, such as the coronavirus disease 2019 (COVID-19) pandemic. A study on Twitter showed that false news spreads more quickly and reaches 100 times as many readers as true news [3]. Therefore, many researchers are working to solve this problem, ranging from analyzing the types of fake news [4] to trying to find the most effective method for detection [5].

In the Arab regions, different rumors and false information have been spread during the COVID-19 pandemic [6]. Rumors and false information have also proliferated in politics, as was the case during the Arab Spring and the Syrian crisis [7]. The Arabic language is considered one of the most commonly spoken languages in the world, and an essential language for Muslims worldwide, who numbered about 1.8 billion in 2015 and are expected to increase to three billion in 2060, according to the Pew Research Center [8]. The Arabic language presents various challenges. First, there are different dialects used in different Arab countries, as well as from regions in the same country. The language also has a rich vocabulary that leads to the occurrence of misleading information from different dialects, making it more difficult for the system to detect [9]. In social media, the complexity of understanding the Arabic language is increased because users on social media, such as Twitter, use two forms of the language: Modern Standard Arabic (MSA), which is the formal language, and Dialect Arabic (DA), which is informal, used in daily communication between users, and is more common than MSA [10].

This paper working to focus in two aspects: features extraction techniques and the datasets used. This is because while classification algorithms play a critical role in fake news detection, their effectiveness heavily relies on the quality and relevance of the features provided to them [11]. Feature extraction techniques enable the models to capture the subtle nuances and contextual cues that differentiate fake news from real news [12]. By incorporating these features, classification algorithms can leverage the rich information present in the data, leading to improved accuracy and robustness in fake news detection. Majority of studies in Arabic fake news

detection have predominantly focused on the application of specific features extractions techniques to identify the veracity of news with no concern on their limitation[13][14][15]. However, this study surpasses mere technical application by shedding light on the strengths and weaknesses of each technique.

On the other hand, focusing on identifying and analyzing available datasets is crucial for advancing research in fake news detection. It allows researchers to improve data quality, benchmark models, and collaborate effectively to address the growing challenge of misinformation in our digital world. There is a need for research which reviews the available dataset in Arabic for fake news detection. For that, this research coming to fill this gap by analyzing available public Arabic datasets according to their domain, size, labels, annotation method, source and the features used. We believe this work provide researchers in the field with valuable insights into the available dimensions to initiate their research endeavors. The contributions of this review paper are as follows:

*1) Investigate* research in Arabic fake news detection by selecting studies that cover different features for detecting fake news and utilize publicly available datasets.

*2) Provide* a clear insight into the available helpful techniques used for feature extraction in detecting fake news, in general, and in the Arabic language, in particular;

*3) Investigate* the publicly available Arabic datasets and show their properties, including a list of the available Arabic fact-checking websites that help build datasets;

*4) Explore* the limitations and prospective avenues concerning datasets, utilized features, and classification methods associated with detecting Arabic fake news.

The rest of this paper is arranged as follows. Section II describes a brief background of fake news, providing its definitions, and impacts. Section III identifies the different approaches for feature extraction. Section IV investigates the works conducted in the field of Arabic fake news detection. Subsequently, Section V provides a list of the available Arabic fact-checking websites used. Section VI investigates the available Arabic dataset in fake news detection. Finally, Section VII discusses some of the limitations and future direction, followed by the conclusion in Section VIII.

## II. BACKGROUND

### A. Fake News Definition

"Fake news" does not have a consistent definition. Scholars define the term differently based on the purpose of their research. For example, one study [16] defined fake news as "a news article that is intentionally and verifiably false." In 2018, a European Commission report defined "fake news" as "all forms of false, inaccurate, or misleading information designed, presented, and promoted to cause public harm intentionally or for profit" [17]. Research in [18] has provided and defined a set of terms related to fake news, including hoax, rumor, spam, misinformation, disinformation, clickbait, satire, propaganda, and hyperpartisan. *Hoaxes* are facts that are either erroneous or inaccurate and are presented as actual facts in news stories [19]. They consist of fabricated information that has been purposefully created to appear true, and because they involve highly intricate and large-scale fabrications, they frequently cause substantial material harm to the victim [20]. *Rumor* is information initiated by a potentially untrustworthy person and circulated among users before it can be verified to be true, false, or unconfirmed [21]. *Spam* refers to any irrelevant or unsolicited messages sent by individuals or groups on social media, such as advertisements, malicious links, or content of poor quality [22]. *Misinformation* is false information that is spread unintentionally, for example, by mistake or updating specific knowledge without intentionally misleading [23]. *Disinformation* is a type of false information spread among users with the main goal of misleading others for some purpose, such as deception of some person [19] or promoting a biased agenda [24]. The concepts of misinformation and disinformation often confuse readers. While both refer to inaccurate or fake information, disinformation is created with malicious intent, whereas this is not necessarily the case with misinformation [25][19]. *Clickbait* is a story whose title or news headline is different from the content itself; it is mainly employed to attract users to access a specific website to increase traffic and, consequently, boost revenue [26][24]. *Satire* is a noteworthy literary tool employed in creating news articles, serving the purpose of both criticism and amusement for readers [27]. This form of discourse often incorporates a substantial amount of irony [28]. *Propaganda* encompasses a form of persuasive communication that employs unidirectional messaging to influence the attitudes, emotions, perspectives, and behaviors of specific target audiences, driven by ideological, political, and religious motives [29][30]. *Hyperpartisan* refers to heavily biased or one-sided narratives [31], particularly in the political arena, denoting strong partiality toward individuals, parties, circumstances, or events [32]. This review will use fake news and false news as a big umbrella for these concepts and will not consider each concept distinctly because there is an overlap in the techniques performed.

### B. Impact of Spreading Fake News

Fake news influences users, leading them to accept biased or false beliefs, changing how they react to true news stories [12]. Spreading fake news not only affects individuals but also harms society. In 2016, fake news drew global attention during the US election [33], when a large amount of fake news was shared on Twitter. The spread of fake news goes beyond the political sector and also applies to other sectors. During the COVID-19 pandemic, a significant amount of fake news and rumors spread throughout the health sector. Approximately 80% of consumers in the US reported that they had read fake news during the pandemic [34]. In addition, during the fire disasters in Australia, some fake maps and pictures were shared on social media [35]. While this may have increased awareness of the disaster, this fake information may also have cost people their lives [35].

In the Middle East, Arabic countries have also been influenced by some fake news, for example, during the COVID-19 pandemic. According to a study on COVID-19 misinformation in Jordan, different Arabic media outlets promoted conspiracy theories regarding the pandemic, with

the most prevalent ideas focusing on the virus's origins [36]. Most of the fake news spreading in Arab countries are related to the political sector, such as news related to the Arab Spring in 2010, Islamic State terrorist group (ISIS) campaigns, and the Beirut explosion of 2020 [37].

### III. FEATURES USED IN FAKE NEWS DETECTION

Feature extraction and selection are techniques used in text mining that have shown effectiveness in enhancing performance in different tasks [10]. Having large features requires more powerful computation and enough memory; hence, there is a vital need to choose appropriate features. In fake news detection, researchers have recently applied feature extraction and selection with good results. In detecting fake news, three approaches are commonly used by researchers: knowledge-based (fact-checking), content approaches, and context approaches [38] [12]. This section will describe each feature in more detail for better understanding. Fig. 1 summarizes these features.



Fig. 1.    Fake news detection approaches.

#### A. Knowledge-based Approaches

This approach works by attempting to discover whether there are facts that support the claim made in a news item [39]. In this approach, fact-checking strategies can be included, which entail finding documents and web pages that support a news item based on information retrieval methods [40]. Few works have used the web to search and retrieve evidence on search query formalization [40]. Study [40] performed a study combining social media conversations with evidence from external sources. The research developed dataset which contain social media conversations and external evidence related to each rumor. Their results showed that combining evidence with rumor is more effective than using rumors or evidence alone. Also, in [41] a system worked by incorporating external evidence from the web with some signals from language style.

#### B. Content-based Approaches

Content-based features focus on information that can be directly extracted from the text, such as linguistic features. The content features can be classified into (textual) linguistic-based [ 74] and visual-based [12].

*1) Linguistic-based features:* Linguistic-based features extract content from the text at different levels, such as characters, words, sentences, and documents [12]. In fake news detection, the common linguistic features used are lexical, syntactic, and semantic [42][12]. Lexical features, considered

as the actual wording (usage) in the text, can be at a word level, such as the total number of words, word length, and the frequency of unique words, or at a character level, such as characters per word [12][43]. One of the commonly used features is bag-of-words (BOW), n-gram, which is used to find the most prominent word contents or expressions [44]. On the other hand, syntactic features include sentence-level features, such as the order of words in a sentence, grammar, and syntactic structure of a sentence [12]. The other features related to linguistics are semantic features, which include using the Natural Language Processing (NLP) technique to extract information, such as opinion mining, and sentiment analysis to extract emotions and opinions from texts [38]. Some research also extract topics in texts or posts on social media using Latent Dirichlet Allocation (LDA) [43]. Word embedding, an example of the distribution semantic technique, is useful in detecting fake news with machine learning and deep learning [45]. The linguistic feature is an effective indicator in detecting fake news during the initial propagation phase with the user features [46]. This method is simple but limited because it requires deep knowledge of the domain; thus, it is difficult to generalize [47]. Another limitation when dealing with linguistic features is that the features extracted from social media, such as Twitter, may not be enough for machine learning approaches because the text is short [48]. Moreover, it cannot be used for detecting fake news that only has images or videos [48].

*2) Visual features:* Visual content features extract features from visual elements, such as videos or photos, using deep learning techniques [49]. Fake news creators use individual vulnerabilities and images or any kind of visual cues to provoke an emotional response from users [12]. To extract visual features from images that can be used in detecting fake news, applications such as clustering scores and similarity distribution histograms are used [49] [12]. In addition, statistical features can be used in fake news verification, such as the multi-image ratio, number of images in an event, image ratio, and long image ratio [49] [12]. Other research has used visual features, such as the polarity of the image and the probability of the image being manipulated [50]. This kind of feature can be integrated with other features, such as text-based features and user-based features, which provide a good result [51]. Few researchers in fake news detection have investigated visual features due to the lack of availability of datasets containing images and video [12] [52]. In addition, the techniques for maintaining these features are more complex than those needed for other features[12] [50].

#### C. Social Context-based Features

Social context-based features are related to the information surrounding news, such as the user's characteristics, the reactions of other users to posts or news, and social network propagation features [42][12].

*1) User-based Features:* User-based features focus on the news publisher to evaluate the credibility of the source [12].

The features are used to confirm credibility and reliability for an individual user using demographics, registration age, number of tweets written by the user, number of followers or following, number of tweets authorized by the user, user photo, user sentiment, tweet repetition, and so on [43]. Most social media platforms provide users with a verification status after they provide information to the social media company to confirm their identity. This verification status is used by researchers in detecting a user's reliability [53]. When using user-based features, it is important to understand that the availability of this information is a critical concern due to privacy and access constraints on some platforms such as Twitter.

*2) Post-based features:* Post-based features can be used to identify the veracity of news from different aspects relevant to social media posts[12]. This includes analyses of user feedback, reactions, opinions, and responses, in general, as indicators of fake news. Some of these features include comments, likes, tagging, user ratings, sentiment, and emotional reactions [42]. Some research dealing with post-based features have dealt with a number of retweets, likes, shares, and others [54][55]. Another unique feature represents the social responses to a post by users who interact with the news story. These features rely on the wisdom of the crowd in detecting fake news and show its effectiveness in detection. In this situation, a few studies have analyzed responses from different perspectives, such as sentiment, emotion, and stance toward news items [12]. In fact, [56] used users' response information as core input data in detecting fake news. They believed that users' responses have rich information that researchers can benefit from.

*3) Network-based features:* Social media users construct networks grounded in relationships, interests, and subjects [12]. Detecting fake news necessitates extracting network-based features to unveil and represent discernible patterns suitable for identification. Different types of networks can be constructed. Shu et al. [12] categorized them into stance network, where the nodes represent tweets related to tweets and the edges represent weights of stance similarities [57] [58]; co-occurrence network, which counts users' written posts related to the same news article [47]; and friendship network, which is based on the following and followers of the users who posted in relation to tweets [59]. After building networks, some matrices can be used for feature representation, such as degree and clustering coefficients, to represent the diffusion network[60]. Research using network-based features is limited compared with other features because of the complexity that can emerge when analyzing the patterns [46]. In addition, finding a dataset that contains enough information for a network is difficult patterns [46].

Research in the field uses each feature alone or incorporates both content and contextual features in detecting fake news [47]. These methods are promising regarding effectiveness but are challenging when relying only on one type of feature in automatic fake news detection. Table I summarizes the shortcomings and merits of each feature.

TABLE I. SHORTCOMINGS AND MERITS FOR FAKE NEWS DETECTION APPROACHES

| Approaches | | Strength and weakness |
|---|---|---|
| **Knowledge-based approach** | *External evidence* | (+) Usually Improve the performance[40]. (+) Enhanced interpretability[98]. (-) Need information retrieval techniques[40]. (-) Need strategies for checking the credibility of sources. |
| **Content-based approaches** | *Textual-based* | (+) Simple method. (+) Appropriate for early detection of fake news[46]. (-) Requires deep knowledge of the domain [47]. (-) May not fully convey the message's context well as underlying intent, leading to misunderstandings [18]. |
| | *Visual-based* | (+) Enhanced the performance[51]. (+) May provide contextual information[50]. (-) Limited dataset containing visual content, especially in Arabic language[12] [52]. (-) Need complex techniques to be maintained from professional & more time for analysis[12][50]. |
| **Social context-based features** | *User-based* | (+) Appropriate for early detection of fake news[46]. (-) There are some privacy concerns. |
| | *Post-based* | (+) Appropriate for early detection of fake news[46]. (-) Need more investigation[12]. |
| | *Network-based* | (+) Good for identifying influential users [12]. (-) Provide low performance in the early stage of detection[46] [48]. (-) Complex to analyzed and require enough information about structure and users' connection [46]. (-) Need advanced computational techniques and expertise in network analysis[12]. |

## IV. FAKE NWES DETECTION AND ARABIC LANGUAGE

Fake news detection in Arabic has gained significant attention recently but is still in its infancy. Some studies were conducted using a simple technique, such as [61], which suggested a solution based on using the rule-based model. The dataset was adapted from Almujaiwel[1], where they built a dictionary that contains a set of fake news with keywords for each news. The rule-based model worked by checking the primary and secondary keys in the dictionary, whether they are in a tweet text or not. In the context of traditional machine learning, various researchers have used common supervised machine learning, such as [62]. The study used word frequency, count vector, and Term Frequency - Inverse Document Frequency (TF-IDF) as features. Moreover, two word-embedding methods were used: FastText and Word2Vec. To collect tweets, the research used Infection Disease Ontology. The result after testing on Logistic Regression (LR), Naïve Bayes (NB), and Support vector machine (SVM) classifiers showed the best accuracy, with 84% for the LR model with a count vector. Alkhair et al. [63] performed another study to detect the content of fake news on YouTube. The study collected data related to rumors about three famous people in Arab countries (topic of the dataset), namely, Fifi Abdu, Adel Imam, and Abdelaziz Boutaflika. The study used features related to content: n-grams with TF-IDF.

[1]https://github.com/salmujaiwel

For the classification process, SVM, Multinomial Naïve Bayes (MNB), and Decision Tree (DT) were used. The best accuracy and precision were registered for SVM across these topics. While the textual features showed their effectiveness, other studies have suggested using user features. El Balloul et al. [64] proposed a model called (CAT) for the credibility analysis of Arabic content on Twitter. CAT is built using a combination of features related to content and user. It has 26 content-based features and 22 user-based features. The research constructed a dataset from Twitter in Arabic, which is considered topic-independent. CAT was trained using NB, SVM, and RF. CAT registered a higher weighted average F-measure of 75.8% using RF. Overall, sentiment was found to be a highly crucial feature in defining credibility, especially the negative sentiment, as well as the URL in the author profile linked to their website. Mouty and Gazdar [65] conducted their research using the same datasets and features of [64]. To enhance the accuracy of the classifier, the study developed an algorithm for discovering the similarity between username and display name in a Twitter account and the similarity score between tweets and Google search results. For classification, they used RF, DT, SVM, and NB. The combination of user/content features and the new two-similarity score features registered 78.71% accuracy for the RF model.

Using similar features, Jardaneh et al. [55] conducted another experiment using LR, Adaptive Boosting (AdaBoost), Random Forest (RF), and DT for classification. The dataset used was from [7]. The research confirmed that the sentiment feature has a beneficial effect on the system's accuracy, especially when ensemble-based machine learning algorithms are employed. Taher et al. [66] performed a study that employed a Harris Hawks Optimizer (HHO) for the feature selection approach. Briefly, the researcher used a combination of features related to user profile, content-based, and linguistic features, namely, TF-IDF, BOW n-grams, and Binary Term Frequency (BTF). Eight machine learning algorithms were tested: eXtreme Gradient Boosting (XGB), NB, K-Nearest Neighbor (KNN), linear discriminant analysis (LDA), DT, LR, SVM, and RF. LR was selected with HHO algorithms, which performed well as the wrapper feature selection approach registered a 5% increase compared with [55], where the accuracy was 82%.

Alzanin and Aqil [67] performed another study for detecting rumors using unsupervised and semi-supervised expectation maximization (EM). The dataset used contained 271,000 tweets belonging to 88 non-rumor events and 89 rumor events, and 16 features related to users and content were used. For the classification tasks, the researchers compared their proposed model with the supervised Gaussian Naïve Bayes (GNB) model. The results showed the proposed model outperformed GNB with an accuracy of 78.6%. In the health sector, [68] focused on cancer treatment information disseminated via social media. The research extracted tweets annotated manually by experts in the field. The total number of datasets was 208 tweets. Given the small set of data, the researchers performed over-sampling to enhance the performance of the model. The features used were TF-IDF and n-grams, while the classifiers were SVM, LR, KNN, Bernoulli

Naïve Bayes (BNB), Stochastic Gradient Decent (SGD), and j48. The research also used an ensemble method using RF, AdaBoost, and Bagging. The results confirmed that the over-sampling process enhanced the performance of all models of machine learning. Meanwhile, RF outperformed the others using 4- and 5-grams based on accuracy. In study [69], they also relied on a set of extracted features from user and textual features. To extract the textual features, they used both classic word embedding (word2vec, fastText, and Keras embedding layer) and context-based embedding (Multilingual Arabic Bidirectional Encoder Representations from Transformers (MARBERT) and ARBERT) with deep learning models[70]. Two deep learning schemes were used: Convolutional Neural Network (CNN) and Bidirectional Long Short-Term Memory (BiLSTM). The result showed that MARBERT with CNN provided the best accuracy of 95.6%. Alawadh et al. [71] performed another experiment using machine learning algorithms and Mini-BERT. The research used the dataset available from [72]. First, the research preprocessed the dataset by applying standard text 2 numeric encoding. Subsequently, DT, NB, RF, linear support vector (LSV), and mini-BERT were applied with three separate splits using the holdout validation technique (70/30, 80/20, 90/10). The result showed that mini-BERT exhibited consistent performance among the splits with increasing training data, while the machine learning classifiers showed varying performances across the splits. The highest accuracy registered with mini-BERT was 98.4%.

Numerous fake news was spreading during the COVID-19 pandemic. These encouraged researchers to build datasets covering this domain, and the deep learning approaches provided them with excellent performance. The researchers started using Neural Network (NN) and a set of language models, such as [9]. They performed experiments to detect misinformation spreading via Twitter related to the COVID-19 pandemic. The dataset size was 8,786 collected Arabic tweets. For classification, the study used eight of the traditional machine learning models, which were MNB, SVM, XGB, SGD, and RF. On the other hand, deep learning models were used, which are CNN, Recurrent Neural Network (RNN) and, Convolutional RNN (CRNN). Both experiments with machine learning and deep learning used different representation features, including word frequency and word embedding. The XGB showed the highest accuracy in detecting misinformation in this study. Study in [73] collected datasets about the COVID-19 pandemic from Twitter on the types of fake news and misinformation to detect fake news. The tweets were annotated in two ways: manual and automatic. For feature extraction, they used count vector, word-level TF-IDF, character-level TF-IDF, and n-gram-level TF-IDF. They chose the best classifier based on performance, which was trained in the manually annotated dataset to automatically annotate the remaining unlabeled dataset as false and real. The study used six classifiers, including RF, NB, LR, Multilayer Perceptron (MLP), XGB, and SVM. As a result, LR exhibited the best performance for both manual and automatic annotated datasets, which registered an F1-score of 93.3%. In this study, the stemming and rooting techniques failed to increase the performance of the classifiers.

TABLE II. SUMMARY OF ARABIC FAKE NEWS DETECTION STUDIES AND CORRESPONDING FEATURES AND MODELS USED.' E' EVIDENCE, 'T' TEXTUAL, 'V' VISUAL, 'U' USER, 'P' POST, 'N' NETWORK

| Study | Dataset | Feature type | | | | | | Models | Platform | Feature Extraction Methods | Evaluation results |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | E | T | V | U | P | N | | | | |
| Alotaibi and Alhammad[61] | 2000 tweets. | - | ✓ | - | - | - | - | Rule-based | Twitter | Building dictionary. | Accuracy = 78.1 % Precision = 70% Recall =98% |
| Alsudias & Rayson [62] | 2000 tweets | - | ✓ | - | - | - | - | LR, SVM, NB | Twitter | Word frequency, Count vector and TF-IDF, FastText, Word2Vec. | Accuracy= 84.03% Precision = 81.04% Recall = 80.03% F1-score= 80.5% |
| Alkhair et al[63] | 3434 comments | - | ✓ | - | - | - | - | SVM, DT, MNB | YouTube | N-grams, TF-IDF | Accuracy= 95.35% Precison = 92.77% Recall= 83.12% |
| El Ballouli et al[64] | 9000 tweets. | - | ✓ | - | ✓ | ✓ | - | NB, SVM, RF | Twitter | - | Precision = 76.1% Recall = 76.3%, F1-score = 75.8% |
| Mouty and Gazdar [65] | Dataset[64] | ✓ | ✓ | - | ✓ | ✓ | - | DT, SVM, NB, RF. | Twitter | - | Accuracy= 78.71% Precision = 78.5% Recall= 78.7%, F1-score = 78.5% |
| Jardaneh et al. [55] | Dataset[7] | - | ✓ | - | ✓ | ✓ | - | RF, DT, LR, AdaBoost | Twitter | - | Accuracy= 76% Precision= 79% Recall= 82% F1-score= 80% |
| Taher et al. [66] | Dataset[7] | - | ✓ | - | ✓ | ✓ | - | XGB, NB, KNN, LDA, DT, LR, SVM, RF. | Twitter | TF, TF-IDF, BTF, N-grams | Accuracy= 82% Precision = 82% Recall= 86%, F1-score= 84% |
| Alzanin and Aqil[67] | 177 events. | - | ✓ | - | ✓ | ✓ | - | EM | Twitter | - | Accuracy=78.6% precision=79.8% recall=80.2% F1-score=78.6% |
| Saeed et al [68] | 208 tweets | - | ✓ | - | - | - | - | SVM, LR, KNN, BNB, SGD, j48, Bagging, AdaBoost, RF | Twitter | TF-IDF, N-gram | Accuracy = 83.50% Precision = 86% Recall= 83% F1-score= 83% |
| Alyoubi et al.[69] | 5000 tweets. | - | ✓ | - | ✓ | - | - | CNN, BiLSTM | Twitter | Keras Embedding Layer, word2vec, FastText, ARBERT. MARBERT, (word and sentence-level) | Accuracy = 95.6% Precision = 95.6% Recall = 95.6% F1-score = 95.6% |
| Alawadh et al.[71] | Dataset[72] | - | ✓ | - | - | - | - | DT, NB, LSV, RF, Mini-BERT | Articles | BERT | Accuracy = 98.43% Precision = 100% Recall =97.5% F1-score = 98.73% |
| Alqurashi et.al[9] | 8786 tweets | - | ✓ | - | - | - | - | MNB, SVM, XGB, SGD, RF. CNN, CRNN, RNN | Twitter | TF-IDF (world-level, N-grams), FastText, word2vec, | Accuracy= 86.2% Precision =67% Recall= 25% F1-score = 37% |
| Mahlous and Al-Laith [73] | 37029 tweets | - | ✓ | - | - | - | - | RF, NB, LR, MLP, XGB, SVM. | Twitter | Count vector, Word-level TF-IDF, N-gram-level TF-IDF. Char-level TF-IDF | **Manual dataset:** Precision=87.8% Recall=87.7% F1= 87.8% **Automatic dataset:** Precision= 93.4% Recall= 93.3% F1-score= 93.3% |
| Elhadad et al.[74] | COVID-19-FAKES [74] | - | ✓ | - | - | - | - | KNN, DT, LR, MNB, LSVM, BNB, Perceptron, NN, XGB, ERF, BME, GB, AdaBoost | Twitter | TF, TF-IDF, N-gram, char-level, word embedding | - |
| Haouari et al [75] | ArCOV19-Rumors [75] | - | ✓ | - | ✓ | ✓ | ✓ | MARBERT, AraBERT, Bi-GCN, RNN+CNN | Twitter | - | Accuracy= 75.7% macro-F1=74% |
| Ameur and Aliane [76] | AraCovid19-MFH [76] | - | ✓ | - | - | - | - | AraBERT, mBERT, Distilbert Multilingual, arabert | Twitter | AraBERT, mBERT, Distilbert Multilingual, arabert Cov19 , mbert | F-score= 95.78% |

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Cov19 , mbert Cov19 | | Cov19 | | |
| Nassif et.al[77] | Kaggle [78] ,10000 records | - | ✓ | - | - | - | - | AraBERT, QaribBert-base, Araelectra, MARBERT. Arabic-Bert, ARBERT, RobertBase, GigaBert-base. | Articles & tweets | - | Accuracy= 98.5% Precison=99.1% Recall=98.2% F1-score=98.6% |
| Touahri and Mazroui [79] | 200 claims, 3380 evidences | ✓ | ✓ | - | - | - | - | Scoring function | Tweets & articles | - | Accuracy= 92.7%. Precision= 54.66%. Recall = 56.13% F1-score= 55.2% |
| Elaziz et al.[80]. | ArCOV19-rumors [75] OSACT4[81] Dataset [73] Dataset [82]. | - | ✓ | - | - | - | - | AraBERT | Twitter & articles | MTL and AraBERT. | Accuracy=95.5% Precision = 96.2% Recall = 96.29% F1-score = 96.28% |
| Amoudi et al[83] | ArCOV19-Rumors [75] | - | ✓ | - | - | - | - | SVM, NB, KNN, DT, RF, SGD, LR, XGB, GRU, RNN, LSTM, Bi-RNN, B-GRU, Bi-LSTM | Twitter | TF-IDF, N-gram, AraVec | Accuracy= 80% Precision= 80% Recall= 72 F1-score= 75% |
| Al-Yahya et al. [45] | ArCOV19-Rumors [75] AraNews[84] ANS [85] COVID-19-Fakes [74] | - | ✓ | - | - | - | - | CNN, GRU, RNN, AraBERT, ArElectra, QARiB, ARBERT, MARBERT | Articles & Twitter | Word level, char-level, Word2Vec, Glove, FastText, doc2vec | Accuracy = 97.5% Precision = 95.6% Recall = 95.6% F1= 95.3% |
| Fouad et .al[86] | Dataset[87] 1980 tweets | - | ✓ | - | - | - | - | LSVC, SVC, MNB, BNB, SGD , DT, RF CNN, LSTM, CNN+LSTM, BiLSTM, CNN + BiLSTM | Articles & Twitter | Word embedding, N-gram | Accuracy= 83.92% |
| Shishah[88] | Covid-19-Fakes [74] ANS [85] Satirical[89] AraNews[84] | - | ✓ | - | - | - | - | BERT with joint learning | Articles & Twitter | - | Accuracy = 85% Precision = 86% Recall = 86% F1-score = 85% |
| Bsoul et al.[90] | 2652 news records | - | ✓ | - | - | - | - | LR, RF, NB, SGD, NN, DT. | Twitter | TF-IDF, BOW. | Precision=84% Recall=78% F1-score= 81% |
| Saadany et al.[89] | Satirical Fake News[89] | - | ✓ | - | - | - | - | MNB, XGB, CNN | News articles | Count vector, TF-IDF, N-grams, char-level, word-level, FastText | Accuracy =98.59% Precision = 98.49% Recall = 98.61% F1-score = 98.49% |
| Khouja[85] | ANS[85] | - | ✓ | - | - | - | - | BERT, LSTM | News title | Word-level, charr-level | Precision = 64.1% Recall = 64.6% F1-score = 64.3% |
| Nagudi et al.[84] | AraNews [84] ATB [2] ANS[85] | - | ✓ | - | - | - | - | AraBERT, mBERT, XLM-R Base, XLM-R Larg | News articles | AraBERT, mBERT, XLM-RBase,, XLM-RLarg | Accuracy = 74.12% F1-score = 70.06% |
| Himdi et al[91] | 1098 records | - | ✓ | - | - | - | - | SVM, NB, RF | Article | POS, Emotion, Polarity, linguistics (syntactic, semantic) | Precision = 79% Recall= 79% F1-score= 79 % |
| Albalawi et al.[52] | 4025 tweets. | - | ✓ | ✓ | - | - | - | AraBERT (different version) ARBERT, MARBER, MARBERTv2, QARIB, Arabic Bert, Arabert Covid-19, mbert Covid-19 Ara-DialectBERT, AraT5, VGG-19, ResNet50 | Twitter | AraBERT (different version), ARBERT, MARBER, MARBERTv2, QARIB, Arabic Bert, Arabert Covid-19 mbert Covid-19, Ara-DialectBERT, AraT5,VGG-19, ResNet50 | Accuracy = 89.8% Precision = 89.47% Recall = 89.87% F1-score= 89.64% |

---

[2]https://www.ldc.upenn.edu/collaborations/past-projects

Study in [74] collected and annotated a set of data from Twitter in Arabic/English language. The research was conducted in two phases. First, a binary classification model was trained on a set of collected ground-truth data, which were obtained from the official websites and the official Twitter accounts of the United Nations (UN), United Nations International Children's Emergency Fund (UNICEF), and World Health Organization (WHO). The second phase carried out the annotation for the, unlabeled tweets. The system used extraction features, such as TF, TF-IDF (n-gram, character level), and word embedding with 13 machine learning algorithms, including KNN, DT, LR, MNB, Linear Support Vector Machines (LSVM), BNB, Perceptron, NN, XGB, Ensemble Random Forest (ERF), Bagging Meta-Estimator (BME), Gradient Boosting (GB), and AdaBoost. Meanwhile, Haouari et al. [75] built the ArCOV19-Rumors dataset, which covered some claims about COVID-19. The research presented a benchmark on claim-level verification and tweet-level verification to exploit the content, user profiles, propagation structure, and temporal features. The Bidirectional Graph CN (Bi-GCN) and RNN+CNN, as well as the AraBERT and MARBERT models were tested on the dataset. The final results showed the best accuracy for AraBERT and MARBERT, with 73% and 75.7%, respectively.

Furthermore, [76] used a pre-trained transformer AraBERT, Multilingual BERT (mBERT), and Distilbert Multilingual under baseline transformer models. The study fine-tuned mBERT and AraBERT on [92]'s datasets containing dialects. The output from this process were two models called AraBERT COV19 and mBERT COV19, which were used in addition to the above-mentioned three models for the experiment (five models overall). The experiment found that the pretrained COVID-19 models were helpful after fine-tuning in detecting false information. Nassif et al. [77] used eight BERT transformer-base models. The research conducted two experiments: using the dataset from Kaggle[3], which was written in English and translated to Arabic using Google translator, and using another dataset, sourced from Twitter and newspaper agency websites. GigaBert-base and QARiB Bert-base provided the best result on the translated dataset and the collected dataset, respectively. Study in [79]focused on the task of claim verification, which involves a collection of claims and corresponding evidence in the form of text snippets sourced from web pages [93]. The researchers achieved an F1-score of 55.2% by employing a scoring function that evaluates the negation and concordance between the claims and their associated text snippets. The determination of negation and concordance levels was performed through a manual list-based approach. Some research focused on feature extraction and selection, such as [80], which used three main methods, including multi-task learning, a transformer-based model, and Fire Hawk Optimization (FHO) algorithm. Three datasets related to COVID-19 and the detection of fake news were used, which include ArCOV19-Rumors, as well as datasets from [73] and [82]. AraBERT was used for feature extraction via multi-tasking learning and fine-tuning approaches, a novel metaheuristic algorithm to select the most pertinent features

from the contextual feature representations. The average accuracy in binary classification using these datasets was 91%. Amoudi et al. [83] performed a comparative study in rumor detection using different machine learning and deep learning models. The dataset used in this research was ArCOV19-Rumors. They used features such as n-grams, TF-IDF, and AraVec word embedding. The first experiment used SVM, KNN, NB, DT, RF, SGD, LR, and XGB for machine learning and evaluated the application of ensemble learning. In addition, six common deep learning models were used: Gated Recurrent Unit (GRU), RNN, LSTM, Bidirectional RNN (Bi-RNN), and Bi-GRU, BiLSTM with seven optimizers. The research showed that over-sampling did not enhance the performance of either the traditional or the deep learning models, and ensemble learning performed better than the single models. LSTM and Bi-LSTM with Root Mean Square Propagation (RMSprop) optimizer provided the best accuracy among the other deep learning models.

Al-Yahya et al. [45] performed another comparison study between the use of the NN model and the transformer-based language model for detecting Arabic fake news. The datasets used in this study included ArCOV19-Rumors, AraNews [84], Arabic News Stance (ANS) corpus [85], and COVID-19-Fakes [74]. The study used a linear model at word and character levels with Glove, Word2Vec, FastText, and document level. For the classification tasks, deep learning models CNN, GRU, and RNN were used. Moreover, QARiB, AraBERT, ArELectra, ARBERT, and MARBERT were used from the transformer-based models. The results showed that the transformer-based models performed better than the NN-based solutions, registering a 95% F1-score for QARiB. Fouad et al. [86] conducted another experiment for detecting fake news. The research used two datasets. The first contains news and tweets collected manually by the researchers and annotated to rumor or non-rumor, while the second dataset was from [87]. The two datasets were merged and used to create a third one. Word embedding and TensorFlow were used for text representation, and a set of traditional machine learning was used, which included Linear Support Vector Classifier (LSVC), SVC, MNB, BNB, SGD, DT, and RF. On the other hand, they examined a set of deep learning models, namely, CNN, LSTM, CNN+LSTM, BiLSTM, and CNN + BiLSTM. As a result, they found that no single model performed optimally over the three categories of datasets from the traditional machine learning models. For deep learning, the BiLSTM model provided the highest accuracy across all three datasets. Meanwhile, Shishah [88] performed another study, proposing a model called JointBERT for detecting the Arabic language. JointBERT in this research used Named Entity Recognition (NER) and Relative Features Classification (RFC) as parameters. The datasets used in this research were COVID-19-Fakes [74], ANS, Satirical [89], and AraNews. The results showed that JointBert outperformed the baseline results. The use of NER increased the performance because of its ability to extract news entities, which supports the model's performance in detecting fake news.

Some researchers have focused on specific types of misinformation, such as Bsoul et al. [90]. They built a dataset for clickbait detection, which facilitated automatic

---

[3] https://www.kaggle.com/c/fake-news/data?select=test.csv

classification and detection of news headlines. The study used BOW and TF-IDF and features related to headlines, such as headline length and with demonstrative pronouns, question words, and question mark. The researchers performed an experiment using SVM, LR, DT, NN. NB, SGD, and RF. These models produced a Macro F1-score of up to 0.81, which shows the effectiveness of using these seven machine learning models in the detection of clickbait news headlines. Research[89] focused in another type of fake news which is satire news. The research used dataset collected from different news websites and working to exploit textual features for the purpose of identifying satire news. For feature extraction they used count vectors, word-level TF-IDF, N-grams, Char-level, with MNB and XGB machine learning. In addition, the research used CNN with pre-trained word embedding which provided best accuracy registered which are 98.59%. The research found that satire news in Arabic language incline to have subjective tone with more positive and negative key terms. Research [85] release a dataset for claim verification, which are derived from subset of news title from Arabic News Text (ANT) corpus[94]. The authors modified news title to generate fake claims. For classification process, they used BERT, and LSTM for training and testing datasets. The study reported that BERT provides best F1-score registered which are 64.3%. Compared to this research, Nagoudi[84] used the same dataset from [85] in addition to their dataset which is automatically generated fake news from real one. The impact of this generated data on verification fake news are tested using transformer-based pre-trained models and compared with human created fake news dataset. The research reported that generated news are positively affect the fake news detection, and achieved better performance than [85] with 70% for F1-score.

In addition, [91] collected a dataset containing articles as fake and real, covering a single topic, which is Al Hajj. For the real articles, they collected articles in three dialects from Arab countries: Saudi, Egyptian, and Jordanian. The veracity of these articles was assessed using different fact-checking platforms. On the other hand, the study used crowdsourcing to create a fake news article based on a real one. The research did not provide a new classification method but suggested a set of features to provide an accurate Arabic classifier. The researchers built a lexical wordlist and an Arabic natural language processing (ANLP) architectural tool to extract the textual features, including POS, emotion, polarity, and syntactic. Based on these features, RF, SVM, and NB were used and tested for each single feature and a combination of features. The best result was registered for RF, with a combination of POS, syntactic and semantic roles, and contextual polarity features, which achieved an F1-score of 79%. Moreover, the research tested against human performance, providing 86% of articles classified correctly. Insufficient emphasis is being placed on the utilization of visual features in the realm of fake news detection, but recently, Albalawi et al. [52] proposed a model based on textual and image features. Their model consisted of three sub-models. For extracting the textual features, the BERT model was used. Meanwhile, two ensembles of pre-trained vision models were used for extracting the visual features (VGG-19 and ResNet50). The final model is a multimodal

model used for concatenating the extracted textual and image features to represent a rumor vector. As a result, their proposed multimodal did not outperform the model with textual-based features. This shows that textual features are still considered pioneering in detecting rumors. Table II provides a brief summary of Arabic fake news detection studies.

## V. ARABIC FACT-CHECKING WEBSITES

When researchers build datasets for fake news detection, they mostly rely on fact-checking websites. Fact-checking websites can be defined as platforms that evaluate the veracity of claims and information spread over social media networks, web articles, and public statements [95]. These kinds of websites usually rely on human experts who check the veracity of the news. Some fact-checking websites employ techniques that rely on evidence-based analysis or assess the credibility of the statements and how they align with the factual reality [95]. In the Arab countries, there are different fact-checking websites that researchers use as a starting point when building datasets for detecting fake news and rumors. It is worthy to provide researchers in the field with a list that can help them during their dataset-building phase. The common Arabic fact-checking websites will be described as follows:

- Norumors[4]: This is a standalone project that started in 2012 as a Twitter account searching for rumors and detecting their veracity. Thereafter, a website was established, which contains archives for different rumors and fake news that spread in Arab countries in general, especially in Saudi Arabia. The aim of this site is to spread awareness about rumors and expose the disseminators of falsehoods.

- Fatabyyano[5]: This is a standalone platform in the field of news fact-checking. First established in 2014 as a single page on Facebook. Fatabyyano uses the Facebook rating system, which has nine rating options, namely, false, partly false, true, false headline, satire, not eligible, opinion, prank generator, and not rated. Fatabyyano is certified by the International Fact-Checking Network (IFCN).

- Misbar[6]: Considered one of the leading fact-checking platforms in the Arab regions, it covers the Middle East and South African countries. The news classification system in Misbar using these classes includes fake, misleading, true, myth, selective, commotion, and satire.

- Akeed [7] : This is a fact-checking platform that is responsible for tracking the credibility of the Jordanian media. It was established with the support of the King Abdullah Fund for Development. The news are rated as false, biased, misleading, ambiguous, incomplete, inciting news, or contains an error.

---

[4] http://norumors.net/?post_type=rumors
[5] https://fatabyyano.net/
[6] https://misbar.com/
[7] https://akeed.jo/public/

- Verify [8] : This is considered the first Syrian fact-checking platform established in 2011 during the Syrian Revolution. In this platform, news are classified into three main classes depending on risk [red (high risk), orange (medium risk), and yellow (low risk)].

- Falso[9]: This is a Libyan platform for monitoring hate speech and a fact-checker in media news. It fact-checks both traditional news from TV and newspaper and also covers social media platforms and other websites.

- FactuelAFP Arabic[10]: This is the Arabic division of the French news agency dedicated to providing news and information. This division holds certification from IFCN and collaborates closely with Facebook to combat misinformation.

- Maharat-news fact-o-meter [11] : This is another fact-checking website certified by IFCN. Their main focus is detecting rumors on online media in Lebanon and other Arab regions. The rating system is based on three labels: true, partially true, and not true.

- Kashif[12]*:* it is an independent Palestinian fact-checking website that main goal to combat misleading information in Palestinian media. The rating system is based on nine labels: Manipulated, incorrect linking, outdated, sarcasm or parody, fabricated, false context, Impersonated and inaccurate content.

- Dabegad [13] : This is a fact-checking website in the Egyptian dialect that started in 2013 and aims to find and expose hoaxes in social media in Middle East countries.

## VI. ARABIC FAKE NEWS DATASETS

It is clear that the first essential step to building an effective fake news detection system is constructing an appropriate dataset. There is still no agreed benchmark dataset for fake news detection in the Arabic language. Researchers in fake news detection usually focus on analyzing the text features only; however, this may no longer be enough, especially with the evolution of social media platforms and the diverse representation of fake news and their complexity. As such, different studies have recently used features related to news publishers, as well as social context features, such as user information and network information [75]. This was evident from the previous sections, which showed promising effectiveness. Accordingly, this section will briefly mention a set of publicly available datasets in the Arabic language. According to type, datasets can be classified into *social media posts* and *news articles*. In social media post datasets, in addition to the text content of posts, user and network information are utilized to detect fake news. Meanwhile, in news article datasets, researchers detect fake news by utilizing the headlines and the body of the articles.

One example of datasets related to posts in social media is [64], which included 9,000 tweets considered topic-independent. The dataset contains information related to the users and content, in addition to the sentiment features. The annotation process was performed manually, classifying the tweets into credible and noncredible. Alzanin and Azmi [67] constructed datasets that consisted of 27,100 posts related to 177 events from different domains, where 88 of these events were considered non-rumors, while 89 were rumors. This dataset also includes some features from users and contents. During the COVID-19 pandemic, researchers were encouraged to construct datasets covering this topic [74][76]. Furthermore, [74] constructed a bilingual dataset that covers the Arabic and English languages. The dataset was collected from Twitter using keywords related to COVID-19. The dataset has been automatically annotated using 13 machine learning classifiers and seven feature extraction techniques, including TF, TF-IDF (n-gram, character level), and word embedding. The dataset has two labels: real and misleading. The tweets' IDs and detection labels are made publicly available to other researchers by the authors.

AraCOVID19-MFH [76] is a multilabel dataset manually annotated to 10 labels. It is an abbreviation of the Arabic COVID-19 Multilabel Fake News and Hate Speech Detection dataset. The dataset consists of 10,828 items of annotated data that include MSA and different dialects. The dataset was collected by searching some keywords from Twitter. This dataset can be used for both hate speech and fake news detection tasks. In [9], another dataset was collected from Twitter by searching for a specific list of keywords. The dataset is labeled as misleading (1311 tweets) and others (7475 tweets). It consists of 8,786 tweets in total, and it is clear that this dataset suffers from imbalance. The researchers only published the ID of the tweets and their corresponding labels (misleading and other). Moreover, in [73], the dataset was collected from Twitter based on the COVID-19 domain. Part of the dataset was annotated manually first and then used to train different machine learning models to automatically produce annotation for the rest of the unlabeled data. They were annotated as fake and genuine and only relied on the text feature. The authors published the tweet's text and its label without the ID. ArCorona [96] is a larger dataset collected manually from Twitter in the health domain during the early stages of COVID-19. The dataset contains 30 million tweets, with 8,000 tweets labeled into 13 labels. The dataset contains different dialects from the Arabic regions.

Alsudais and Rayson [62] also manually collected a dataset that contains one million tweets about COVID-19, among which 2,000 were annotated to 895 false, 0 unrelated, and 789 true. The tweets were collected using keywords related to infectious diseases. All the above-mentioned datasets rely on textual content; however, there is a dataset called ArCOV19-Rumors [75], which consists of 9,414 tweets that were manually labeled as false (1,753), true (1,831), and others (5,830), which are related to 138 claims This dataset works on two levels: claim-level verification and tweet-level verification, both of which have two labels, namely, true and false. and their correspond relevant tweets. Meanwhile, tweet-level verification contains the tweet and the propagation

---

8 https://verify-sy.com/
9 https://falso.ly/
10 https://factcheck.afp.com/ar/list
11 https://maharat-news.com/fact-o-meter
12 https://kashif.ps/category/facts-in-en/
13 https://dabegad.com/

network (retweets, replies). Moreover, study in [82] built dataset about Covid-19. The datasets are labeled for binary classification and multiclass, based on seven questions. The datasets available in four language which is English, Arabic, Dutch, and Bulgarian. The main focus for this paper is fake news, so 4966 Arabic tweets are classified into (815) contains fake information and (2602) not contains false information.

TABLE III. SUMMARY OF SOCIAL MEDIA DATASETS FOR ARABIC FAKE NEWS DETECTION. 'E' EVIDENCE, 'T' TEXTUAL, 'V' VISUAL, 'U' USER, 'P' POST, 'N' NETWORK, 'M' MANUALLY, 'A' AUTOMATICALLY

| Dataset | Domain | Size | Annotated | Annotation type | | Labels | Features used | | | | | | Source |
| | | | | M | A | | E | T | V | U | P | N | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Dataset [64] | Multi-domain | 9K tweets | 9K tweets | ✓ | - | (5400) credible (3600) noncredible | - | ✓ | - | ✓ | ✓ | - | Twitter |
| Dataset [67] | Multi-domain | 271K tweets 177 events | 271K tweets | ✓ | - | (88) events non-rumor (89) events rumor | - | ✓ | - | ✓ | ✓ | - | |
| COVID-19-FAKE [74] | Health (covid-19) | 220K tweets | 220K tweets | - | ✓ | Misleading Real | - | ✓ | - | - | - | - | |
| AraCOVID19-MFH[76] | Health (covid-19) | 300K tweets | 10828 tweets | ✓ | - | 10 different labels | - | ✓ | - | - | - | - | |
| Dataset [9] | Health (covid-19) | 4.5M tweets | 8.8K tweets | ✓ | - | (1,311) misleading, (7,475) other | - | ✓ | - | - | - | - | |
| Dataset [73] | Health (covid-19) | 36066 tweets | 36066 tweets | ✓ | ✓ | (20417) Fake (15649) Not fake | - | ✓ | - | - | - | - | |
| ArCorona[96] | Health (covid-19) | 30M tweets | 8K tweets | ✓ | - | 13 different labels | - | ✓ | - | - | - | - | |
| Dataset [62] | Health (covid-19) | 1M tweets | 2K tweets | ✓ | - | (316) False (895) True (789) Unrelated | - | ✓ | - | - | - | - | |
| ArCovid19-Rumors[75] | Health (covid-19) | 1M tweets 138 claims | 9414 tweets | ✓ | - | (1753) false (1831) true (5830) other | - | ✓ | - | ✓ | ✓ | ✓ | |
| Dataset [82] | Health (covid-19) | 4966 tweets | 4966 tweets | ✓ | - | (2609) No (815) Yes | - | ✓ | - | - | - | - | |
| Dataset [7] | Politic | 3358 tweets | 2708 tweets | ✓ | - | (1570) credible (1138) non-credible | - | ✓ | - | - | - | - | |

TABLE IV. SUMMARY OF NEWS ARTICLE DATASETS FOR ARABIC FAKE NEWS DETECTION. 'E' EVIDENCE, 'T' TEXTUAL, 'V' VISUAL, 'S' STANCE

| Dataset | Domain | Size | Annotation type | Labels | Features | | | | Source |
| | | | | | E | T | V | S | |
|---|---|---|---|---|---|---|---|---|---|
| AraNews [84] | Multi-domain | 5187957 news articles | Manually | False, True | - | ✓ | - | - | 50 newspapers. |
| Dataset[72] | Multi-domain | 323 articles | | (100) reliable (223) unreliable | - | ✓ | - | - | Kashif fact-checking website, social media, WhatsApp group, and news site. |
| AraFact [97] | Multi-domain | 6222 claims | | (4037) false (1891) partly-false (198) True (90) Sarcasm (6) unverifiable | ✓ | ✓ | ✓ | - | 5 Arabic fact-checking websites |
| Dataset [98] | Politic (Syria War) | 422 claims 3042documents | | (219) false claim. (203) true claim Documents: (1239) false (1803) True | - | ✓ | - | ✓ | 2 websites VERIFY for false claim and REUTERS[14] for True claim |
| ANS [85] | Multi-domain | 4547 claims (3786) pairs (claim, evidence) | | (3072) True (1475) False | - | ✓ | - | ✓ | News headlines from media sources in Middle east and ANT corpus. |
| AraStance [99] | Multi-domain | 910 claims 4063 pair (claim, articles) | | (606) False claim (304) True claims Articles: (2421) False (1642) True | - | ✓ | - | ✓ | 3 fact-checking websites Aranews, Dabegad, Norumors, REUTERS |
| Satirical fake News[89] | Politic | 3185 articles | | (3185) fake. (3710) real. | - | ✓ | - | - | 2 satirical news websites for fake news. And Official news site for real. |
| AFND [100] | Multi-domain | 606912 articles | | (207310) credible (167233) not credible (232369) undecided | - | ✓ | - | - | 134 Arabic online news sources |
| Dataset. [7] | Multi-domain | 175 blog posts | | (100) credible (57) fairly credible (18) non-credible | - | ✓ | - | - | Arabic news articles |

---

[14] http://ara.reuters.com

For the news article dataset, AraNews [84] is a large dataset collected from different countries, covering different topics. The dataset relied on a list of 50 newspapers from 15 Arab countries, the United Kingdom, and the United States. Based on this list, 5,187,957 news articles were collected and labeled as true or false. Also, research [72] published dataset which contains 323 articles from different domains. The dataset was labeled manually from two experts into reliable and unreliable. Different sources are used for collecting this dataset such as Kashif fact-checking websites, social media, and news sites. Ali et al. [97] produced a dataset called AraFacts that contains claims collected from five Arabic fact-checking websites. The dataset comprises URLs for fact-checking articles and links to evidence pages sourced from various outlets, enabling the extraction of images and supporting evidence. This making them the first to collect datasets in Arabic that combined texts, image contents, and evidence. Research in [52] used AraFacts to extract images and text. The dataset contained 6,222 claims annotated by professional fact-checkers manually.

Typically, researchers focus independently on specific tasks such as stance detection, fact-checking, document retrieval, and source credibility, but study [98] making these tasks available to be integrated using unified corpus. The study created a corpus consisted of 442 claims that covered topics related to the Syrian war and some political issues in the Middle East. Each claim was labeled as false or true based on factuality. In addition, 2,042 articles retrieved for these claims were annotated based on their stance into agree, disagree, discuss, and unrelated. In addition, the dataset included a rationale attribute that enabled the fact-checking system to provide explanations for its decisions. This attribute encompassed the display of extracted sentences or phrases from the retrieved documents, which served as illustrative examples of the detected stance. Similarly, ANS [85] is a dataset of news titles that were paraphrased and altered. ANS covers several topics collected from online news sites. The dataset contains two versions. Based on claims verification, 4,547 records were labeled fake and not fake; the other version consisted of claims and evidence pairs containing 3,786 records. The difference between ANS and Baly et al. [98] is that ANS generated fake and true claims from true news. In the same vein, the AraStance dataset [99] contained 4,063 pairs of claims and articles from multiple countries, covering topics in politics, health, sports, and others. The annotation process was performed manually according to the veracity of the claims, whether they are true or false, and according to the article's stance on the claims, which ended with four labels: agree, disagree, discuss, and unrelated. Satirical News [89] is a hand-built dataset that includes fake articles. The fake news articles were collected from two satirical news websites: Al-Hudood [15] and Al-Ahram Al-Mexici[16]. For the real news dataset, the research used an open-source datasets[17]. Arabic Fake News Dataset (AFND) [100] is a collection of news stories in Arabic that are available to the public and were gathered from Arabic news websites. A dataset used for detecting article credibility, it consisted of 606,912 articles labeled as credible (207,310), not credible (167,233), and undecided (232,369). The researchers used the Misbar fact-checking platform for classifying the articles into these three classes. Some research built datasets that contained both social media posts and news articles, such as Al Zaatari et al.'s dataset [7], which consisted of blog posts and tweets related to the Syrian crisis. This kind of corpora's main goal is to analyze the credibility of the news. It consisted of 2,708 tweets and 175 blog posts; the datasets, in general, are labeled as credible and not credible. Tables III and IV provide a summary of the social media posts and news article datasets. In the context of news article datasets, the range of available features is comparatively narrower compared to those accessible for social media datasets. Thus, our focus was confined to textual, visual, stance, and evidence features for the purpose of comparison.

## VII. DISCUSSION

Based on the previous sections, fake news detection in the Arabic language is still in its nascent stages compared with other languages. Although a multitude of efforts have been exerted in Arabic fake news detection, the process still suffers from various limitations. We can categorize them into limitations related to datasets, feature extraction, and classification algorithms based on the literature.

### A. Datasets

There is no benchmark dataset in the Arabic language, and most of the datasets used in previous studies are not available online. Researchers need to enrich the fake news detection field by making their data available using any platform, such as their own page in GitHub[18] or MASADER[19] repository. Publicly available datasets enable other studies to exert robust efforts, and their results can be compared with others using the same datasets. This is one of the important factors to measure the improvement of performance among different studies. Moreover, the processes for collecting and preparing datasets are not always mentioned clearly in the research papers. Researchers in this field need accurate guidelines to undertake this process, for example, by providing a list of common sources for extracting appropriate news, the annotation process, the cleaning and preprocessing phase, and so on. In addition, most of the datasets suffer from unbalanced classes, where the real news category is usually larger than the false one such as [9] [64]. Moreover, when screening parts of these datasets, there were numerous instances for the same news, especially those collected from social media platforms such as Twitter. This repetition in the values of the news texts can be attributed to the fact that the same news can propagate among such platforms, and because the collection process usually depends on the prepared list of keywords, the probability of having the same news with the same exact text is higher. In some situations, having duplicated contents will decrease the performance of the classifier [101]. Another problem related to the datasets is the domain they cover. For example, Tables III and IV show that most of the datasets focused on covering the COVID-19 pandemic, while the others focused on whether

---

the data was a post from social media or articles covering different domains with unbalanced categories for each. These researchers argued that they adapted this process to find common features across various domains. However, they ignored the fact that each sector has its own features, which can help in detecting fake news related to them. Having datasets that cover a specific domain is important. Specialized terminology is prevalent within specific topic domains, and possessing technical expertise is essential for discerning the authenticity of each news item [18]. Hence, it is crucial to construct datasets tailored to the unique requirements of a particular topic domain, such as the SCIFACT dataset [102].

Moreover, Arabic countries cover different regions with different dialects; thus, datasets covering MSA and dialects require a complex process to mitigate this problem. We still need an Arabic dataset that covers the MSA language and others that cover different dialects. The dataset publishers usually release only the IDs for the tweets in accordance with Twitter policy. When other researchers try to extract the same tweets, they are already deleted or protected by their owners. This situation decreases the number of available datasets and is one of the immense problems that need to be solved by applying a repository that does not contain critical features and saves only appropriate and valuable ones. Working to extract more records to balance the dataset affected by this deletion is time-consuming. The creation of an unaltered dataset for the purpose of detecting fake news on social media platforms represents a crucial milestone in establishing a benchmark dataset. This endeavor aids in effectively assessing the performance of models and their ability to verify the accuracy of information with the collaboration of social media platforms.

*B. Feature Extraction and Classification Process*

Most previous studies have focused on extracting textual features in articles and social media posts, which is not usually sufficient to detect fake news [18]. Rather, it is important to combine features such as those related to users, posts, networks for social posts, and the metadata related to the articles, including the publisher's name, URLs, headlines, body, and comments. Based on the literature review, only a few researches are using a combination of these features. In addition, when looking for the visual contents, images provide an improvement in the performance of the classifier in other languages, which need more investigation in Arabic [50]. There are a few research in Arabic with this attribute, such as [52]. Extracting features from the replies of users is also an important aspect [56]. The only publicly available Arabic dataset that has considered users' responses in their dataset is ArCOVID19-Rumors. Employing sentiment analysis, stance detection, and emotion recognition with fake news detection still needs more investigation, especially in relation to the responses of the crowd. Moreover, features related to networks are still not investigated in fake news in the Arabic language, except for the work performed by [75]. Traditional machine learning approaches, such as SVM, LR, and KNN, are the most used in Arabic research. Based on the review, most research in Arabic has used language models, such as AraBERT, MARBERT, mBERT, and QARiB. Research using the ensemble techniques is scarce, and there is a need for more

studies applying this model. Applying the ensemble methods using traditional machine learning and deep learning is another window of opportunity that needs to be investigated for fake news detection because ensemble methods have the ability to improve prediction performance with regard to some characteristics, such as overfitting avoidance, computational advantages, and representation [103].

## VIII. CONCLUSION

This review revealed that a few studies related to detecting fake news have been conducted in the Arabic language, indicating that Arab researchers should allocate more attention to this issue. Most Arabic studies have focused on social media platforms, particularly Twitter. Most Arabic researchers have investigated events related to politics (e.g., Syrian crisis) and health (e.g., COVID-19 pandemic) and have created their own dataset for testing the proposed models. This may not be an effective approach because it results in many different datasets with a range of accuracy measurement results. A benchmark dataset that contains as many features as possible to help detect fake news should be used. Our future objectives encompass the development of an Arabic dataset encompassing diverse extraction features, including visual attributes, enabling us to explore their influence when combined with other features for the purpose of detecting fake news on social media.

## REFERENCES

[1] Kepios, "Global Social Media Statistics," Datareportal, 2022. https://datareportal.com/social-media-users (accessed Nov. 10, 2023).

[2] M. Walker and K. E. Matsa, "News Consumption Across Social Media in 2021," Pew Research Center, 2021. https://www.pewresearch.org/journalism/2021/09/20/news-consumption-across-social-media-in-2021/ (accessed Feb. 11, 2023).

[3] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," Science (80-. )., vol. 359, no. 6380, pp. 1146–1151, 2018, doi: 10.1126/science.aap9559.

[4] A. Al-Rawi, A. Fakida, and K. Grounds, "Investigation of COVID-19 Misinformation in Arabic on Twitter: Content Analysis," JMIR Infodemiology, vol. 2, no. 2, p. e37007, Jul. 2022, doi: 10.2196/37007.

[5] S. Kula and R. K. P. W. M. Choraś Michałand Kozik, "Sentiment Analysis for Fake News Detection by Means of Neural Networks," in Computational Science -- ICCS 2020, 2020, pp. 653–666.

[6] P. Nakov, F. Alam, S. Shaar, G. Martino, and Y. Zhang, "A Second Pandemic? Analysis of Fake News About COVID-19 Vaccines in Qatar," ArXiv, vol. abs/2109.1, 2021.

[7] A. Al Zaatari et al., "Arabic Corpora for Credibility Analysis," in Proceedings of the Tenth International Conference on Language Resources and Evaluation ({LREC}'16), May 2016, pp. 4396–4401, [Online]. Available: https://aclanthology.org/L16-1696.

[8] M. LIPKA and C. HACKETT, "Why Muslims are the world's fastest-growing religious group," Pew Research Center, 2017. https://www.pewresearch.org/short-reads/2017/04/06/why-muslims-are-the-worlds-fastest-growing-religious-group/ (accessed Nov. 10, 2023).

[9] S. Alqurashi, B. Hamoui, A. S. Alashaikh, A. Alhindi, and E. A. Alanazi, "Eating Garlic Prevents COVID-19 Infection: Detecting Misinformation on the Arabic Content of Twitter," ArXiv, vol. abs/2101.0, 2021.

[10] H. ALSaif and T. Alotaibi, "Arabic Text Classification using Feature-Reduction Techniques for Detecting Violence on Social Media," Int. J. Adv. Comput. Sci. Appl., vol. 10, May 2019, doi: 10.14569/IJACSA.2019.0100409.

[11] S. K. Hamed, M. J. Ab Aziz, and M. R. Yaakub, "A review of fake news detection approaches: A critical analysis of relevant studies and

highlighting key challenges associated with the dataset, feature representation, and data fusion.," Heliyon, vol. 9, no. 10, p. e20382, Oct. 2023, doi: 10.1016/j.heliyon.2023.e20382.

[12] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake News Detection on Social Media: A Data Mining Perspective," SIGKDD Explor. Newsl., vol. 19, no. 1, pp. 22–36, Sep. 2017, doi: 10.1145/3137597.3137600.

[13] R. Mouty and A. Gazdar, "Survey on Steps of Truth Detection on Arabic Tweets," in 2018 21st Saudi Computer Society National Computer Conference (NCC), 2018, pp. 1–6, doi: 10.1109/NCG.2018.8593060.

[14] S. Althabiti, M. Alsalka, and E. Atwell, "A Survey: Datasets and Methods for Arabic Fake News Detection," Int. J. Islam. Appl. Comput. Sci. Technol., vol. 11, pp. 19–28, 2023.

[15] R. A. M. San Ahmed, "A Novel Taxonomy for Arabic Fake News Datasets," Int. J. Comput. Digit. Syst., vol. 14, no. 1, pp. 159–166, 2023, doi: 10.12785/ijcds/140115.

[16] Z. I. Mahid, S. Manickam, and S. Karuppayah, "Fake News on Social Media: Brief Review on Detection Techniques," in 2018 Fourth International Conference on Advances in Computing, Communication Automation (ICACCA), 2018, pp. 1–5, doi: 10.1109/ICACCAF.2018.8776689.

[17] A. D'Ulizia, M. C. Caschera, F. Ferri, and P. Grifoni, "Fake news detection: a survey of evaluation datasets," PeerJ. Comput. Sci., vol. 7, pp. e518–e518, Jun. 2021, doi: 10.7717/peerj-cs.518.

[18] T. Murayama, "Dataset of Fake News Detection and Fact Verification: A Survey," arXiv, 2021, [Online]. Available: http://arxiv.org/abs/2111.03299.

[19] S. Kumar, R. West, and J. Leskovec, "Disinformation on the Web: Impact, Characteristics, and Detection of Wikipedia Hoaxes," in Proceedings of the 25th International Conference on World Wide Web, 2016, pp. 591–602, doi: 10.1145/2872427.2883085.

[20] V. Rubin, Y. Chen, and N. Conroy, "Deception detection for news: Three types of fakes," Proc. Assoc. Inf. Sci. Technol., vol. 52, pp. 1–4, 2015, doi: 10.1002/pra2.2015.145052010083.

[21] C. Buntain and J. Golbeck, "Automatically Identifying Fake News in Popular Twitter Threads," in 2017 IEEE International Conference on Smart Cloud (SmartCloud), 2017, pp. 208–215, doi: 10.1109/SmartCloud.2017.40.

[22] R. Ghanem and H. Erbay, "Context-dependent model for spam detection on social networks," SN Appl. Sci., vol. 2, no. 9, p. 1587, 2020, doi: 10.1007/s42452-020-03374-x.

[23] P. Hernon, "Disinformation and misinformation through the internet: Findings of an exploratory study," Gov. Inf. Q., vol. 12, no. 2, pp. 133–139, 1995, doi: https://doi.org/10.1016/0740-624X(95)90052-7.

[24] S. Volkova, K. Shaffer, J. Y. Jang, and N. Hodas, "Separating Facts from Fiction: Linguistic Models to Classify Suspicious and Trusted News Posts on Twitter," in Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Jul. 2017, pp. 647–653, doi: 10.18653/v1/P17-2102.

[25] L. Wu, F. Morstatter, K. M. Carley, and H. Liu, "Misinformation in Social Media: Definition, Manipulation, and Detection," SIGKDD Explor. Newsl., vol. 21, no. 2, pp. 80–90, Nov. 2019, doi: 10.1145/3373464.3373475.

[26] Y. Chen, N. J. Conroy, and V. L. Rubin, "Misleading Online Content: Recognizing Clickbait as 'False News,'" in Proceedings of the 2015 ACM on Workshop on Multimodal Deception Detection, 2015, pp. 15–19, doi: 10.1145/2823465.2823467.

[27] J. Brummette, M. DiStaso, M. Vafeiadis, and M. Messner, "Read All About It: The Politicization of 'Fake News' on Twitter," Journal. \& Mass Commun. Q., vol. 95, no. 2, pp. 497–517, 2018, doi: 10.1177/1077699018769906.

[28] C. Burfoot and T. Baldwin, "Automatic Satire Detection: Are You Having a Laugh?," in Proceedings of the ACL-IJCNLP 2009 Conference Short Papers, Aug. 2009, pp. 161–164, [Online]. Available: https://aclanthology.org/P09-2041.

[29] C. Lumezanu, N. Feamster, and H. Klein, "#bias: Measuring the Tweeting Behavior of Propagandists," Proc. Int. AAAI Conf. Web Soc. Media, vol. 6, no. 1, pp. 210–217, 2021, [Online]. Available: https://ojs.aaai.org/index.php/ICWSM/article/view/14247.

[30] G. S.Jowett and V. O'Donnell, Propaganda & persuasion. SAGE, 2014.

[31] M. Potthast, J. Kiesel, K. Reinartz, J. Bevendorff, and B. Stein, "A Stylometric Inquiry into Hyperpartisan and Fake News," in Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Jul. 2018, pp. 231–240, doi: 10.18653/v1/P18-1022.

[32] S. Zannettou, M. Sirivianos, J. Blackburn, and N. Kourtellis, "The Web of False Information: Rumors, Fake News, Hoaxes, Clickbait, and Various Other Shenanigans," J. Data Inf. Qual., vol. 11, no. 3, May 2019, doi: 10.1145/3309699.

[33] Amy Watson, "Fake news worldwide - Statistics & Facts," 2020. https://www.statista.com/topics/6341/fake-news-worldwide/#dossierKeyfigures (accessed Nov. 10, 2021).

[34] A. Waston, "Fake news in the U.S. - statistics & facts | Statista," Jun. 16, 2021. https://www.statista.com/topics/3251/fake-news/ (accessed Nov. 10, 2021).

[35] G. Rannard, "Australia fires : Misleading maps and pictures go viral," BBC Trending, 2020. https://www.bbc.com/news/blogs-trending-51020564.

[36] M. Sallam et al., "High Rates of COVID-19 Vaccine Hesitancy and Its Association with Conspiracy Beliefs: A Study in Jordan and Kuwait among Other Arab Countries," Vaccines, vol. 9, 2021.

[37] J. Klausen, "Tweeting the Jihad: Social Media Networks of Western Foreign Fighters in Syria and Iraq," Stud. Confl. Terror., vol. 38, no. 1, pp. 1–22, Jan. 2015, doi: 10.1080/1057610X.2014.974948.

[38] M. A. Alonso, D. Vilares, C. Gómez-Rodríguez, and J. Vilares, "Sentiment Analysis for Fake News Detection," Electronics, vol. 10, p. 1348, 2021.

[39] J. Z. Pan, S. Pavlova, C. Li, N. Li, Y. Li, and J. Liu, "Content Based Fake News Detection Using Knowledge Graphs BT - The Semantic Web – ISWC 2018," 2018, pp. 669–683.

[40] J. Dougrez-Lewis, E. Kochkina, M. Arana-Catania, M. Liakata, and Y. He, "PHEMEPlus: Enriching Social Media Rumour Verification with External Evidence," in Proceedings of the Fifth Fact Extraction and VERification Workshop (FEVER), May 2022, pp. 49–58, doi: 10.18653/v1/2022.fever-1.6.

[41] K. Popat, S. Mukherjee, A. Yates, and G. Weikum, "DeClarE: Debunking Fake News and False Claims using Evidence-Aware Deep Learning," in Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 2018, pp. 22–32, doi: 10.18653/v1/D18-1003.

[42] A. Bondielli and F. Marcelloni, "A Survey on Fake News and Rumour Detection Techniques," Inf. Sci., vol. 497, no. C, pp. 38–55, Sep. 2019, doi: 10.1016/j.ins.2019.05.035.

[43] C. Castillo, M. Mendoza, and B. Poblete, "Information Credibility on Twitter," in Proceedings of the 20th International Conference on World Wide Web, 2011, pp. 675–684, doi: 10.1145/1963405.1963500.

[44] H. Ahmed, I. Traoré, and S. Saad, "Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques," 2017.

[45] M. Al-Yahya, H. Al-Khalifa, H. Al-Baity, D. AlSaeed, and A. Essam, "Arabic Fake News Detection: Comparative Study of Neural Networks and Transformer-Based Approaches," Complexity, vol. 2021, p. 5516945, 2021, doi: 10.1155/2021/5516945.

[46] S. Kwon, M. Cha, and K. Jung, "Rumor Detection over Varying Time Windows," PLoS One, vol. 12, no. 1, pp. 1–19, 2017, doi: 10.1371/journal.pone.0168344.

[47] N. Ruchansky, S. Seo, and Y. Liu, "CSI: A Hybrid Deep Model for Fake News Detection," Proc. 2017 ACM Conf. Inf. Knowl. Manag., 2017.

[48] Y. Liu and Y.-F. Wu, "Early Detection of Fake News on Social Media Through Propagation Path Classification with Recurrent and Convolutional Networks," 2018.

[49] Z. Jin, J. Cao, Y. Zhang, J. Zhou, and Q. Tian, "Novel Visual and Statistical Image Features for Microblogs News Verification," IEEE Trans. Multimed., vol. 19, no. 3, pp. 598–608, 2017, doi: 10.1109/TMM.2016.2617078.

[50] S. K. Uppada, P. Patel, and S. B., "An image and text-based multimodal model for detecting fake news in OSN's," J. Intell. Inf. Syst., 2022, doi: 10.1007/s10844-022-00764-y.

[51] Y. Wang et al., "EANN: Event Adversarial Neural Networks for Multi-Modal Fake News Detection," in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery &amp; Data Mining, 2018, pp. 849–857, doi: 10.1145/3219819.3219903.

[52] R. M. Albalawi, A. T. Jamal, A. O. Khadidos, and A. M. Alhothali, "Multimodal Arabic Rumors Detection," IEEE Access, vol. 11, pp. 9716–9730, 2023, doi: 10.1109/ACCESS.2023.3240373.

[53] Z. S. Ali, A. Al-Ali, and T. Elsayed, "Detecting Users Prone to Spread Fake News on Arabic Twitter," in Proceedinsg of the 5th Workshop on Open-Source Arabic Corpora and Processing Tools with Shared Tasks on Qur'an QA and Fine-Grained Hate Speech Detection, 2022, pp. 12–22.

[54] R. Alghamdi and O. Alrwais, "Towards Automatic Rumor Detection in Arabic Tweets," Int. J. Data Min. Manag. Syst., vol. 1, no. 5, pp. 1–14, 2022.

[55] G. Jardaneh, H. Abdelhaq, M. Buzz, and D. Johnson, "Classifying Arabic tweets based on credibility using content and user features," in 2019 IEEE Jordan International Joint Conference on Electrical Engineering and Information Technology, JEEIT 2019 - Proceedings, 2019, pp. 596–601, doi: 10.1109/JEEIT.2019.8717386.

[56] H. Kidu, H. Misgna, T. Li, and Z. Yang, "User Response-Based Fake News Detection on Social Media BT  - Applied Informatics," 2021, pp. 173–187.

[57] E. Tacchini, G. Ballarin, M. L. Della Vedova, S. Moret, and L. de Alfaro, "Some Like it Hoax: Automated Fake News Detection in Social Networks," ArXiv, vol. abs/1704.0, 2017.

[58] M. Davoudi, M. R. Moosavi, and M. H. Sadreddini, "DSS: A hybrid deep model for fake news detection using propagation tree and stance network," Expert Syst. Appl., vol. 198, p. 116635, 2022, doi: https://doi.org/10.1016/j.eswa.2022.116635.

[59] N. Zhong, G. Zhou, W. Ding, and J. Zhang, "A Rumor Detection Method Based on Multimodal Feature Fusion by a Joining Aggregation Structure," Electronics, vol. 11, no. 19, 2022, doi: 10.3390/electronics11193200.

[60] S. Kwon, M. Cha, K. Jung, W. Chen, and Y. Wang, "Prominent Features of Rumor Propagation in Online Social Media," in 2013 IEEE 13th International Conference on Data Mining, 2013, pp. 1103–1108, doi: 10.1109/ICDM.2013.61.

[61] F. L. Alotaibi and M. M. Alhammad, "Using a Rule-based Model to Detect Arabic Fake News Propagation during Covid-19," Int. J. Adv. Comput. Sci. Appl., vol. 13, no. 1, 2022, doi: 10.14569/IJACSA.2022.0130114.

[62] L. Alsudias and P. Rayson, "COVID-19 and Arabic Twitter: How can Arab World Governments and Public Health Organizations Learn from Social Media?," Jul. 2020, [Online]. Available: https://aclanthology.org/2020.nlpcovid19-acl.16.

[63] M. Alkhair, K. Meftouh, K. Smaïli, and N. Othman, "An Arabic Corpus of Fake News: Collection, Analysis and Classification," in Arabic Language Processing: From Theory to Practice, 2019, pp. 292–302.

[64] R. El Ballouli, W. El-Hajj, A. Ghandour, S. Elbassuoni, H. Hajj, and K. Shaban, "CAT: Credibility Analysis of Arabic Content on Twitter," in WANLP 2017, co-located with EACL 2017 - 3rd Arabic Natural Language Processing Workshop, Proceedings of the Workshop, Apr. 2017, pp. 62–71, doi: 10.18653/v1/w17-1308.

[65] R. Mouty and A. Gazdar, "Employing the Google Search and Google Translate to Increase the Performance of the Credibility Detection in Arabic Tweets BT  - Computational Collective Intelligence," 2022, pp. 781–788.

[66] T. Thaher, M. Saheb, H. Turabieh, and H. Chantar, "Intelligent Detection of False Information in Arabic Tweets Utilizing Hybrid Harris Hawks Based Feature Selection and Machine Learning Models," Symmetry (Basel)., vol. 13, no. 4, 2021, doi: 10.3390/sym13040556.

[67] S. M. Alzanin and A. M. Azmi, "Rumor Detection in Arabic Tweets Using Semi-Supervised and Unsupervised Expectation–Maximization," Know.-Based Syst., vol. 185, no. C, Dec. 2019, doi: 10.1016/j.knosys.2019.104945.

[68] F. Saeed, W. M.S., M. Al-Sarem, and E. Abdullah, "Detecting Health-Related Rumors on Twitter using Machine Learning Methods," Int. J.

Adv. Comput. Sci. Appl., vol. 11, Jan. 2020, doi: 10.14569/IJACSA.2020.0110842.

[69] S. Alyoubi, M. Kalkatawi, and F. Abukhodair, "The Detection of Fake News in Arabic Tweets Using Deep Learning," Appl. Sci., vol. 13, no. 14, 2023, doi: 10.3390/app13148209.

[70] M. Abdul-Mageed, A. Elmadany, and E. M. B. Nagoudi, "ARBERT & MARBERT: Deep Bidirectional Transformers for Arabic," ArXiv, vol. abs/2101.0, 2020.

[71] H. M. Alawadh, A. Alabrah, T. Meraj, and H. T. Rauf, "Attention-Enriched Mini-BERT Fake News Analyzer Using the Arabic Language," Futur. Internet, vol. 15, no. 2, 2023, doi: 10.3390/fi15020044.

[72] R. Assaf and M. Saheb, "Dataset for Arabic Fake News," in 2021 IEEE 15th International Conference on Application of Information and Communication Technologies (AICT), 2021, pp. 1–4, doi: 10.1109/AICT52784.2021.9620228.

[73] A. Mahlous and A. Al-Laith, "Fake News Detection in Arabic Tweets during the COVID-19 Pandemic," Int. J. Adv. Comput. Sci. Appl., vol. 12, 2021, doi: 10.14569/IJACSA.2021.0120691.

[74] M. K. Elhadad, K. F. Li, and F. Gebali, "COVID-19-FAKES: A Twitter (Arabic/English) Dataset for Detecting Misleading Information on COVID-19," 2020.

[75] F. Haouari, M. Hasanain, R. Suwaileh, and T. Elsayed, "ArCOV19-Rumors: Arabic COVID-19 Twitter Dataset for Misinformation Detection," in Proceedings of the Sixth Arabic Natural Language Processing Workshop, Apr. 2021, pp. 72–81, [Online]. Available: https://aclanthology.org/2021.wanlp-1.8.

[76] M. S. Hadj Ameur and H. Aliane, "AraCOVID19-MFH: Arabic COVID-19 Multi-label Fake News & Hate Speech Detection Dataset," Procedia Comput. Sci., vol. 189, pp. 232–241, 2021, doi: https://doi.org/10.1016/j.procs.2021.05.086.

[77] A. B. Nassif, A. Elnagar, O. Elgendy, and Y. Afadar, "Arabic fake news detection based on deep contextualized embedding models," Neural Comput. Appl., vol. 34, no. 18, pp. 16019–16032, 2022, doi: 10.1007/s00521-022-07206-4.

[78] "Fake News | Kaggle," 2018. https://www.kaggle.com/c/fake-news/data?select=test.csv (accessed Feb. 06, 2023).

[79] I. Touahri and A. Mazroui, "EvolutionTeam at CLEF2020-CheckThat! lab: Integration of Linguistic and Sentimental Features in a Fake News Detection Approach.," 2020.

[80] M. Abd Elaziz, A. Dahou, D. A. Orabi, S. Alshathri, E. M. Soliman, and A. A. Ewees, "A Hybrid Multitask Learning Framework with a Fire Hawk Optimizer for Arabic Fake News Detection," Mathematics, vol. 11, no. 2, 2023, doi: 10.3390/math11020258.

[81] F. Husain, "OSACT4 Shared Task on Offensive Language Detection: Intensive Preprocessing-Based Approach," in Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection, May 2020, pp. 53–60, [Online]. Available: https://aclanthology.org/2020.osact-1.8.

[82] F. Alam et al., "Fighting the COVID-19 Infodemic: Modeling the Perspective of Journalists, Fact-Checkers, Social Media Platforms, Policy Makers, and the Society," in Findings of the Association for Computational Linguistics: EMNLP 2021, Nov. 2021, pp. 611–649, doi: 10.18653/v1/2021.findings-emnlp.56.

[83] G. Amoudi, R. Albalawi, F. Baothman, A. Jamal, H. Alghamdi, and A. Alhothali, "Arabic rumor detection: A comparative study," Alexandria Eng. J., vol. 61, no. 12, pp. 12511–12523, 2022, doi: https://doi.org/10.1016/j.aej.2022.05.029.

[84] E. M. B. Nagoudi, A. A. Elmadany, M. Abdul-Mageed, T. Alhindi, and H. Cavusoglu, "Machine Generation and Detection of Arabic Manipulated and Fake News," CoRR, vol. abs/2011.0, 2020, [Online]. Available: https://arxiv.org/abs/2011.03092.

[85] J. Khouja, "Stance Prediction and Claim Verification: An Arabic Perspective," in Proceedings of the Third Workshop on Fact Extraction and VERification (FEVER), Jul. 2020, pp. 8–17, doi: 10.18653/v1/2020.fever-1.2.

[86] K. M. Fouad, S. F. Sabbeh, and W. Medhat, "Arabic Fake News Detection Using Deep Learning," Comput. Mater. \& Contin., vol. 71, no. 2, pp. 3647–3665, 2022, doi: 10.32604/cmc.2022.021449.

[87]  F. Rangel, P. Rosso, A. Charfi, W. Zaghouani, B. Ghanem, and J. Snchez-Junquera, "Overview of the track on author profiling and deception detection in arabic," 2019.

[88]  W. Shishah, "JointBert for Detecting Arabic Fake News," IEEE Access, vol. 10, pp. 71951–71960, 2022, doi: 10.1109/ACCESS.2022.3185083.

[89]  H. Saadany, C. Orasan, and E. Mohamed, "Fake or Real? A Study of Arabic Satirical Fake News," in Proceedings of the 3rd International Workshop on Rumours and Deception in Social Media (RDSM), Dec. 2020, pp. 70–80, [Online]. Available: https://aclanthology.org/2020.rdsm-1.7.

[90]  M. A. Bsoul, A. Qusef, and S. Abu-Soud, "Building an Optimal Dataset for Arabic Fake News Detection," Procedia Comput. Sci., vol. 201, pp. 665–672, 2022, doi: https://doi.org/10.1016/j.procs.2022.03.088.

[91]  H. Himdi, G. Weir, F. Assiri, and H. Al-Barhamtoshy, "Arabic Fake News Detection Based on Textual Analysis," Arab. J. Sci. Eng., 2022, doi: 10.1007/s13369-021-06449-y.

[92]  S. Alqurashi, A. Alhindi, and E. A. Alanazi, "Large Arabic Twitter Dataset on COVID-19," ArXiv, vol. abs/2004.0, 2020.

[93]  M. Hasanain et al., "Overview of CheckThat! 2020 Arabic: Automatic identification and verification of claims in social media," 2020.

[94]  A. Chouigui, O. Ben Khiroun, and B. Elayeb, "ANT Corpus: An Arabic News Text Collection for Textual Classification," in 2017 IEEE/ACS 14th International Conference on Computer Systems and Applications (AICCSA), 2017, pp. 135–142, doi: 10.1109/AICCSA.2017.22.

[95]  Z. Guo, M. Schlichtkrull, and A. Vlachos, "A Survey on Automated Fact-Checking," Trans. Assoc. Comput. Linguist., vol. 10, pp. 178–206, 2022, doi: 10.1162/tacl_a_00454.

[96]  H. Mubarak and S. Hassan, "ArCorona: Analyzing Arabic Tweets in the Early Days of Coronavirus (COVID-19) Pandemic," in Proceedings of the 12th International Workshop on Health Text Mining and Information Analysis, Apr. 2021, pp. 1–6, [Online]. Available: https://aclanthology.org/2021.louhi-1.1.

[97]  Z. Sheikh Ali, W. Mansour, T. Elsayed, and A. Al - Ali, "AraFacts: The First Large Arabic Dataset of Naturally Occurring Claims," in Proceedings of the Sixth Arabic Natural Language Processing Workshop, Apr. 2021, pp. 231 - 236, [Online]. Available: https://aclanthology.org/2021.wanlp-1.26.

[98]  R. Baly, M. Mohtarami, J. Glass, L. Màrquez, A. Moschitti, and P. Nakov, "Integrating Stance Detection and Fact Checking in a Unified Corpus," in Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), Jun. 2018, pp. 21–27, doi: 10.18653/v1/N18-2004.

[99]  T. Alhindi, A. Alabdulkarim, A. Alshehri, M. Abdul-Mageed, and P. Nakov, "AraStance: A Multi-Country and Multi-Domain Dataset of Arabic Stance Detection for Fact Checking," in Proceedings of the Fourth Workshop on NLP for Internet Freedom: Censorship, Disinformation, and Propaganda, Jun. 2021, pp. 57–65, doi: 10.18653/v1/2021.nlp4if-1.9.

[100] A. Khalil, M. Jarrah, M. Aldwairi, and M. Jaradat, "AFND: Arabic fake news dataset for the detection and classification of articles credibility," Data Br., vol. 42, p. 108141, 2022, doi: https://doi.org/10.1016/j.dib.2022.108141.

[101] H.-Y. Lu, C. Fan, X. Song, and W. Fang, "A novel few-shot learning based multi-modality fusion model for COVID-19 rumor detection from online social media.," PeerJ. Comput. Sci., vol. 7, p. e688, 2021, doi: 10.7717/peerj-cs.688.

[102] D. Wadden et al., "Fact or Fiction: Verifying Scientific Claims," in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), Nov. 2020, pp. 7534–7550, doi: 10.18653/v1/2020.emnlp-main.609.

[103] O. Sagi and L. Rokach, "Ensemble learning: A survey," Wiley Interdiscip. Rev. Data Min. Knowl. Discov., vol. 8, no. 4, pp. 1–18, 2018, doi: 10.1002/widm.1249.

# Enhancing Quality-of-Service in Software-Defined Networks Through the Integration of Firefly-Fruit Fly Optimization and Deep Reinforcement Learning

Mahmoud Aboughaly[1*], Shaikh Abdul Hannan[2]

Mathematics Department-Faculty of Science, Ain Shams University, Cairo, Egypt[1]
Assistant Professor, Department of Computer Science and Information Technology,
Al-Baha University, Al-Baha, Kingdom of Saudi Arabia[2]

*Abstract*—The Software Defined Networking (SDN) paradigm has emerged as a critical tool for meeting the dynamic demands of network management with respect to efficiency and flexibility. Quality of Service (QoS) optimization, which encompasses essential features including bandwidth allocation, latency, and packet loss, is a major problem in SDN systems due to its direct influence on network application performance and user experience. To deal with these important issues, this paper tackles the critical problem of Software-Defined Networks (SDNs) Quality-of-Service (QoS) optimization, which is a critical factor affecting network application performance and user experience. Within the Firefly-Fruit Fly Optimised Deep Reinforcement Learning (DQ-FFO-DRL) framework, a novel combination of optimization techniques derived from Fruit Fly and Firefly behaviors with Deep Q-Learning is presented in this suggested approach, which is called Deep Q-Learning. The framework effectively investigates ideal network configurations by utilizing the distinct advantages of the Fruit Fly and Firefly optimization components, while the Deep Q-Learning component dynamically adjusts to changing network circumstances by drawing conclusions from prior experiences. Extensive testing and modeling reveal that the DQ-FFO-DRL approach performs very well in SDNs compared to conventional QoS management solutions. When it comes to negotiating the always changing world of resource allocation, network usage, and overall network performance, this algorithm demonstrates exceptional adaptability. The suggested system, which is implemented in Python, offers an advanced and flexible method for enhancing QoS in SDN systems.

*Keywords—Software Defined Network (SDN); Quality of Service (QoS); firefly-fruit fly optimization; Deep Reinforcement Learning (DRL); adaptive QoS enhancement; network optimization*

## I. INTRODUCTION

Network administrators have been using traffic engineering approaches to enhance resources management efficiency in order to cope with the ever-increasing volume of network traffic. In communications networks, traffic engineering synchronises the packet forward pathways of various streams within the network to enhance the total level of service provided by network users. Routing optimisation remains a long-term research topic along with one of the main obstacles in network utilities optimisation through traffic engineering [1]. Using conventional routing techniques, each router decides how to forward packets on its own, disregarding the judgments made by other routers. Even though this dispersed routing technique is scalable as it can be used on any size networks, it is challenging to handle network management of resources in an effective and flexible manner and optimise network routing as a whole. Software defined networking (SDN) was initially proposed as an efficient way to handle the whole network by dividing control and information layers in the network [2], that separates information transfer from control functions. These factors make an improved network management models more necessary [3]. SDN logically decouples the network's operational plane and the information plane to give an overall picture of the system and enhances network programming capabilities for network administration and operation. Networking policies can be deployed dynamically and with efficiency using this SDN approach. SDN makes it possible to control forwarding of packets centrally and see the entire network, but creating the best routing scheme is not easy. Limited shortest path issues are how the issue of routing is formulated in many current publications, yet these issues typically have an optimal solution that is NP-hard [4]. Furthermore, while the generic multi-commodity flow issue has a conventional solution that takes the network's operation as a constant model with fluctuating traffic, these models are unable to effectively depict good network function under complicated and variable traffic conditions [5].

In recent years, deep reinforcement learning (DRL), that blends deep neural networks and reinforcement learning (RL), has been used to create traffic engineering strategies [6]. The development of the DRL method offers a fresh approach to optimising extremely complex transportation issues. By enhancing routing policy effectiveness in a model-free and focused on experiences way, DRL-based route methods are able to develop and adjust to complicated networks. Using the DRL approach in an SDN-based networks has shown to significantly improve routing optimisation performance, according to recent studies. It should be highlighted that networks performance loss can happen throughout the learning procedure, especially in the beginning phases, owing to the characteristics of reinforcement learning (RL), that entails experimentation in the manner of identifying the most effective approach [7]. When training an infrastructure, one should not risk network efficiency deterioration when improper routing policies immediately boost packet loss and

end-to-end delay throughout the network itself. This decreases the system's dependability. Specifically, in the event of a system's topology modification, the DRL agent that is running ought to retrain how to optimise routing. Continuous network performance deterioration is caused by the longer time it takes for convergence to occur when the features of internet traffic become more complicated. In addition, using DRL-based route optimisation, which requires investigation, can have severe effects in systems that carry QoS-sensitive information [8].

The need for better Quality of Service (QoS) provisioning in Software-Defined Networks (SDNs) is growing in importance in the ever-changing world of today's networks. As numerous applications develop and information traffic increases, guaranteeing effective and adaptive QoS administration has become a critical task. Conventional methods of maximising quality of service frequently lack the flexibility required to handle the intricate and shifting dynamics of contemporary network settings [9]. The study explores the incorporation of novel Firefly-FruitFly Optimised Deep Q-learning approaches to provide adaptive QoS improvements in SDNs in order to overcome these constraints. By separating both data and control planes, Software-Defined Networking (SDN) completely transformed the way networks are managed and controlled. This has made it possible to have centralised control over the network. Due to this division, network assets can be used with never-before-seen flexibility and control, allowing for dynamic setups and modifications in response to changing needs [10]. But even with SDN's built-in benefits, maintaining excellent QoS is still difficult because of the complex interactions between different network variables as well as the constantly changing nature of internet traffic.

Among the most important metrics for assessing the efficiency of a network is quality of service, which includes a number of factors such as latency, bandwidth, dependability, and safety. Optimising these variables to match particular service level commitments and guarantee the best possible experience for users is necessary for successful QoS management. Conventional QoS administration frequently uses preset or static regulations, which may not be able to adjust to shifting network circumstances. This could result in less-than-ideal resource usage and possible performance issues. Techniques inspired by environment have become more popular in recent decades for resolving challenging optimisation issues [11]. Based by the typical behaviours of fruit flies and fireflies, accordingly, the Firefly and FruitFly Optimisation algorithms have been found to be remarkably effective in solving a wide range of optimisation problems. These methods use the ideas of repellent and attraction to identify the best answers in challenging optimisation scenarios, imitating the typical actions of these bugs. By incorporating such bio-inspired methods into social media, there is a chance to improve the flexibility and effectiveness of QoS control in SDNs. Moreover, the utilisation of RNN-LSTM in the field of networks has demonstrated exceptional capacity to tackle intricate and variable optimisation issues. In DRL, robots are trained to communicate with the surroundings and make successive choices in order to maximise aggregate rewards. The network's controllers can be given the ability to

generate wise and flexible choices depending on the needs and circumstances of the network's infrastructure in actual time by utilising DRL in connection with SDN [12]. The study attempts to establish a new architecture which not just optimises QoS in SDNs but additionally responds to shifting traffic trends and network behaviour by merging DRL into the Firefly-FruitFly Optimisation methods.

Software-Defined Networks (SDNs) are a revolutionary paradigm in computer networking that have arisen to address the increasing needs of dependable and effective network services. The optimisation of Quality of Service (QoS) attributes, which include jitter, latency, throughput and packet loss and have a direct impact on network application performance and user experience, is a key challenge in SDNs. This paper presents a novel framework that combines Deep Q Learning with Firefly-Fruit Fly Optimisation to transform QoS enhancement in SDNs. The bio-inspired optimisation algorithms of Firefly-Fruit Fly Optimisation efficiently explore and identify optimal network configurations, drawing inspiration from the natural behaviours of both fruit flies and fireflies. At the same time, Deep Q-Learning's adaptive learning mechanism keeps learning from previous experiences and network interactions. This allows the framework to make wise decisions in real time and adeptly adjust to the changing network conditions. The combination of these cutting-edge methods offers an SDN environment that is more responsive, effective, and optimised, constituting a major breakthrough in the field of software-defined networking. The study would then go into experimental validation, performance metrics, and comparisons with traditional QoS management strategies in order to show how this new framework performs exceptionally well and how it ultimately improves user experience and network efficiency. The following lists the main findings of the suggested investigation:

- The main contribution of the paper is to advance the field of SDNs by offering a fresh, flexible, and clever framework for enhancing QoS. The combination of deep reinforcement learning and optimisation inspired by nature offers a novel strategy for tackling the problems related to quality of service (QoS) in contemporary network settings.

- The article presents a novel framework that incorporates the nature-inspired optimisation technique known as Firefly-Fruit Fly Optimization. Quality of Service (QoS) management in Software-Defined Networks (SDNs) presents intrinsic issues that could be properly explored and identified through the implementation of this optimization technique.

- The QoS management method gains a dynamic and self-learning mechanism with the integration of Deep Reinforcement Learning (DRL) into the framework. Based on prior interactions and experiences, DRL continuously adjusts to shifting network conditions, offering a clever and flexible method of improving QoS.

- The management of crucial QoS factors, including as packet loss, throughput, jitter, and delay, is the study's

primary objective. The framework combines DRL and Firefly-Fruit Fly Optimisation to dynamically optimise these parameters in order to improve network performance as a whole.

The subsequent sections of the paper are organized to provide a comprehensive exploration and validation of the proposed framework. The research unfolds in a structured manner, with the following key segments: A summary of relevant routing of networks work is given in Section I. The paper's uniqueness and the issues with related works are discussed in Section II, Problem statement in Section III. The SDN QoS improvement mechanism is described in Section IV. Explain the suggested method's effectiveness rating in Section V, and the overall research's conclusions and future work is given in Section VI and Section VII respectively.

## II. RELATED WORKS

A relatively new development in networking for computers is the software-defined network (SDN), which separates forwarding of information from centralised control to provide a very adaptable and controllable networks architecture. A great deal of study has been conducted to provide effective routes and allocate resources for SDNs. To guarantee application-driven QoS regardless of the face of computer hacking, situation-aware networks administration continues to encounter significant obstacles. To deal with this matter, Hossain and Wei [13] Utilise technology related to reinforcement learning (RL) to provide smart networks administration and scenario knowledge from the standpoint of routing control. In the modelling parts, the effectiveness of the suggested RL-enabled route management approach is assessed by taking into account various situations. Despite the efforts to enhance the recommended routing method, one possible limitation is the substantial resource commitment needed to fully evaluate the approach's effectiveness in a sizable testbed. The extent and velocity of the experimental and assessment process may be constrained by the monetary and time commitment.

Conventional networks for routing use limited data to determine routing, that can cause a delay in adapting to network traffic fluctuation and a lack of assistance for apps' QoS needs. Casas-Velasco et al. [14] presents reinforced learning and Software-Defined Networking's Intelligent Routing (RSIR), a revolutionary method for navigating in SDN. In order to generate routing choices, RSIR incorporates an Understanding Planes into SDN and specifies a Reinforcement Learning (RL)-based routing algorithms that considers link-state data. The method computes and installs the best routes previously in the forwarded devices by utilising the artificial intelligence offered by RL, the worldwide view and management of the network that is given by SDN, and its relationship with the surrounding environment. Utilising actual traffic matrix for imitation, RSIR was thoroughly assessed. The findings indicate that when bandwidth available, postpone, and damage are taken into account separately or together for the estimation of optimum pathways, RSIR works better than Dijkstra's method in terms of exertion, link productivity, loss of packets, and time. The outcomes indicate that RSIR is a desirable option for SDN smart routing. The

need for substantial computing power for implementing Deep Reinforcement Learning (DRL) to enhance RSIR choices constitutes a potential limitation for future studies, especially for bigger networks. Additionally, integrating traffic forecasts for selecting a path may add to the method's complexities.

The concept of Software Defined Networking (SDN), which centralises intellect in software-driven controllers in order to increase network adaptability and address various network difficulties, continues to gain traction in both research and the IT sector. SDN is considered one of the driving forces behind 5G networks. The efficiency and utilisation of the network may be improved and optimised with the help of machine learning (ML) technologies. Network administration and operation tricky issues have shown to be a huge cooperative challenge for Neural Networks (NN) and Reinforcement Learning (RL) in specific. Bouzidi et al.[15], an SDN-based principles insertion method that uses Deep Q-Network (DQN) agents to acquire the best routes and redirect congestion in order to increase networks utilisation. The method primarily uses NN to proactively forecast traffic bottlenecks. To achieve this, authors initially outline the connection problem that considers Quality-of-Service (QoS) as a Linear Programme (LP), with the goal of minimising both link utilisation and from beginning to end (E2E) time. Following that, author suggests a effective heuristic approach to resolve it. Emulation-based mathematical results utilising Mininet and ONOS controllers show that the suggested method can greatly enhance network capabilities by reducing lost packets, E2E postponement, and connection utilisation. One possible limitation for further study is the sophisticated optimisation procedure needed to set up a Distributed Deep Q-Network (DQN) agent, which adds complication to the installation process and allows for efficient management of the quantity and location of SDN routers and related information plane switching.

Conventional routing algorithms use only a small amount of data to determine routing, that results in a delayed response to traffic fluctuation and a constrained capacity to fulfil apps' quality of service needs. In order to overcome these drawbacks, researchers presented, a Reinforcement Learning (RL)-based router approach for SDN. Yet, when confronting vast actions and state areas, RL-based systems typically see a rise of training time. Casas-Velasco et al.[16] Presents Deep Reinforcement Learning and Software Defined Networking Intelligent Routing (DRSIR), an alternative routing method. In SDN, DRSIR specifies an algorithm for routing that gets beyond the drawbacks of RL-based systems by utilising Deep RL (DRL). In order to generate innovative, efficient, and smart routing that adjusts to continuous congestion changes, DRSIR takes path-state indicators into account. Emulation was used to assess DRSIR utilising both artificial and actual traffic patterns. According to the outcomes, this method operates better in terms of flex, loss of packets, and latency than the routes that utilise Dijkstra's algorithm and RSIR. Furthermore, the outcomes show that DRSIR offers a workable and realistic approach for networking in SDN. The intricacy and mathematical requirements of expanding DRSIR to accommodate multiple paths routing, multiple levels DRL plans, and the integration of travel type data could be a

hindrance for subsequent research, making it difficult to control the method's learning settings and assess effectiveness across a variety of traffic conditions.

The requirement for quality of services resulting from an exponential rise in network traffic makes routing optimisation increasingly vital. With the latest advancement of software-defined networking (SDN) technologies, networking equipment like switches can now be flexible configuration via programming interfaces, allowing for centralised administration and operation. Kim et al. [17] Provide a routing optimisation on an SDN using deep reinforcement learning (DRL). Under the suggested approach, the agent that handles DRL determines an ideal set of connection weights to strike a compromise for the networking's lost packets and total latency by learning how the overall traffic load of network routers and network efficiency are related. Installing the flow-rules onto the SDN-enabled shifts, an SDN controller uses an array of link strengths to decide how to distribute pathways. They create an M/M/1/K queue-based network framework and use it to execute the DRL training procedure offsite unless it converges in order to circumvent the extremely lengthy learning procedure for DRL in the event of a topologies modification. The outcomes of the experiment show that in a number of network structures, the suggested routing technique works better than both a network demand-based RL algorithms and a traditional hop-count forwarding technique. To guarantee the efficacy and usability of the suggested routing approach in multiple network settings, additional assessment and verification of the approach over a larger range of configurations is necessary. This represents a single negative.

Numerous applications that operate in real time utilise reinforcement learning (RL), a way of learning without supervision. A challenge involving making choices is at the heart of RL. In reinforcement learning, the participant engages with its surroundings continuously and decides what to do future based on past input regarding rewards. Younus et al.[18] RL Software-Defined Wireless Sensor Networks (SDWSNs) optimise their routing pathways through training. They merge SDN and RL, when routing lists are generated by applying RL to the SDN controller. In addition, they suggest four distinct incentive mechanisms to optimise the efficiency of networks. When contrasted with RL-based forwarding algorithms, RL-based SDWSN enhances the network's efficiency by 23% to 30% as a matter of lifetime. Since it's able to effectively learn the network path at the level of the controller, RL-based SDWSN operates well. Furthermore, compared to RL-based WSN, it offers a faster network integration rate. One possible disadvantage of SDWSN is that its centralised control can give rise to problems with scalability and higher communication above you, which could restrict its use in larger and intricate network systems.

The literature reviewed here emphasises how software-defined networking (SDN) and reinforcement learning (RL) are becoming increasingly important in routing optimisation to fulfil the demands of effective resource allocation, quality of service (QoS), and network adaptability. Scholars like Hossain and Wei, Casas-Velasco, Bouzidi et al., and Kim et al. investigate several RL-based methods for smart routing in

SDN, taking traffic patterns, link-state data, and deep reinforcement learning (DRL) into account. When compared to conventional routing methods, these studies show increases in network efficiency, less packet loss, and decreased latency. On the other hand, difficulties like the need for more computer power, the difficulty of optimisation, and issues with scalability for larger networks are recognised. Furthermore, Younus et al. investigate the incorporation of RL in SDN within the framework of wireless sensor networks (SDWSNs), demonstrating improved network performance using RL-based routing choices. Scalability concerns are brought up by the centralised control of SDWSNs, despite their benefits. Overall, the literature highlights the promise for intelligent and adaptive network management through the integration of SDN and RL, while also highlighting the need for more study to overcome obstacles and optimise these strategies for a range of network situations.

## III. PROBLEM STATEMENT

In the context of larger and more complicated network systems, the research recognises a number of noteworthy issues related to the recommended methodologies. One notable drawback is the high resource consumption, which includes processing power, time, and monetary inputs needed to evaluate the efficacy and flexibility of the suggested methods. These resource limitations could prevent the technologies from being widely used, which would limit their scalability. Potential challenges arise from the complexity of implementation and management processes brought about by the integration of Deep Reinforcement Learning (DRL) into routing methods. Moreover, it has been observed that the centralised administration of Software-Defined Networking (SDN) [16] systems leads to increased inefficiencies and scalability problems, which restricts the adaptability of solutions in complex network environments. The study underscores the need for additional evaluation and verification of the suggested techniques in various network contexts to guarantee their efficacy and pragmatic suitability. In spite of these obstacles, the study presents a novel and cutting-edge approach to addressing Quality of Service (QoS) in SDNs by combining Fruit Fly Optimised Q-Learning with Firefly. Using the intelligence of deep reinforcement learning and the flexible abilities of nature-inspired optimisation, this method actively maximises network efficiency based on current traffic requirements. With its ability to adapt to changing traffic patterns and network conditions, the combination that has been developed provides a stable and efficient network architecture that shows promise in addressing the issues raised in the issue's formulation.

## IV. OUTLINE OF THE PROPOSED MECHANISM

The suggested hybrid method improves Quality of Service (QoS) in Software Defined Networks (SDN) by combining Deep Reinforcement Learning (DRL) with Firefly-Fruit Fly Optimisation. The strengths of both optimisation strategies are used in this integrated methodology. The optimisation of network parameters through the combined intelligence of fruit flies and firefly is known as Firefly-Fruit Fly Optimisation. The quality of a potential solution is represented by the brightness of each firefly in a model of firefly interaction

called Firefly Optimisation. This brightness-based method aids in fine-tuning network setups, with enhanced packet loss, throughput, latency, jitter, and brightness directing increases in certain parameters. Thorough Reinforcement Network configuration adaptation is heavily reliant on learning. It makes use of neural networks' capacity to learn and make judgments in a sequential manner that is consistent with QoS objectives. DRL improves performance by allowing the network to dynamically adjust to changing circumstances. It regularly evaluates the network's condition, pinpoints areas in need of modification, and puts new policies into effect as necessary. The combination of DRL and Firefly-Fruit Fly Optimisation is a synergistic method for improving QoS. Firefly Optimisation is excellent at finding viable solutions, and DRL gives you the tools to put those answers into practise and modify them quickly. When combined, they maximise throughput while reducing latency, jitter, and packet loss in order to optimise network parameters. With this integration, network efficiency and flexibility will be maintained, which will ultimately result in a more dependable and seamless user experience and a considerable improvement in QoS in SDN. The suggested technique's workflow is depicted in Fig. 1.



Fig. 1. Workflow of the proposed system.

### A. Data Collection

In this research, an evaluation is conducted to compare the performance of SDN (Software-Defined Networking) and traditional non-SDN network configurations. The examination takes place within the network infrastructure of the engineering faculty at Universitas Muhammadiyah Malang (UMM). The primary goal is to assess the capabilities of both SDN and non-SDN networks when subjected to the same network topology. To facilitate this assessment, the study employs simulation techniques, specifically using the SDN network emulator and the MiniNet Floodlight controller. A

range of Quality of Service (QoS) parameters, including latency delay, jitter, throughput, and packet loss, is employed to gauge the QoS metrics of the latter network [19].

### B. Feature Extraction Based on Hybrid Firefly-Fruit Fly Optimization

A nature-inspired optimization technique called Fruit Fly Optimization (FFO) is based on the swarming behaviours of fruit flies. It is intended to resolve intricate optimization issues by mimicking the motion and communication of fruit flies as they seek for the optimal solution [20].

Step 1: Establish the primary FOA settings and randomly assign the fruit fly swarm's starting position L.

$$u - axis, v - axis$$

Step 2: Give your own fruit fly the ability to go in any direction in search of nourishment by employing Eq. (1) and Eq. (2):

$$u_p = U - axis + RV \qquad (1)$$

$$v_p = V - axis + RV \qquad (2)$$

$$P = 1,2, \dots, h$$

where, $h$ is the magnitude of the fruit fly swarm.

Step 3: Given the inherent uncertainty in determining the exact location of food, we can calculate the distance (represented as $Distance^p$) of the fruit fly from its starting point. This calculation allows us to establish a judgment value for the concentration of the smell (denoted as $F_p$). Let's assume that Si is the reciprocal of $Distance^p$, as follows in Eq. (3) and Eq. (4):

$$Distance^p = \sqrt{u_p^2 + v_p^2} \qquad (3)$$

$$F_p = \frac{1}{Distance^p} \qquad (4)$$

Step 4: By entering the smell concentration judgment value ($F_p$) into the scent concentration judgment function (also known as the Fitness function), one may obtain the scent intensity ($SL_p$) of each unique fruit fly site in Eq. (5).

$$SL_p = Fn\ (F_p) \qquad (5)$$

Step 5: Determine which fruit fly in the swarm has the strongest scent concentration on an individual basis in Eq. (6):

$$[\ Best_{SL}\ Best_{index}] = Maximum\ (SL_p) \qquad (6)$$

Step 6: Preserve the optimal fruit fly's position (u, v) and highest scent intensity level. The swarm then takes off for that destination in Eq. (7):

$$SL(\ Smell)Best = \ Best_{SL} \qquad (7)$$

$$u - axis = u(Best_{index})$$

$$v - axis = v(Best_{index})$$

To reiterate the execution of steps 2 to 6, initiate iterative optimization. The loop concludes when either the number of iterations reaches the maximum allowed limit or when the

current concentration of scent no longer surpasses the concentration obtained in the previous iteration.

*1) Firefly Algorithm (FA):* The Firefly algorithm, created by Xin-She [21], is predicated on the idealised behaviour of firefly flashing qualities. These flashing traits could be summed up as follows, for ease of understanding:

- In the firefly world, gender is not a factor; every firefly is universally drawn to others. This means that a firefly will be attracted to any other firefly, regardless of its gender.

- Attractiveness in this context is directly tied to the luminosity of fireflies. When two fireflies flash, the one with less brightness will be naturally inclined to move closer to the brighter one. This attractiveness factor is directly proportional to their respective brightness levels, and it diminishes as the distance between them increases. In the absence of a firefly brighter than itself, a firefly will resort to random movement.

- The brightness of a firefly is intricately linked to the landscape of the objective function they aim to optimize. In other words, the features of the terrain, such as peaks and valleys, will determine the brightness of a firefly.

To simplify the understanding, it could assume that a firefly's appeal is determined by its brightness or the intensity of its light, which, in turn, relates to the encoded objective function. In the most straightforward scenario for optimization, we can express the brightness (R) of a firefly at a specific position (u) as $R(u) \propto 1/f(u)$. However, it's important to note that attractiveness is a relative concept, and it should be contingent on the distance ($e_{xy}$) between two fireflies, $x$ and $y$. Similar to how light intensity diminishes as you move away from its source and is affected by the medium it travels through; it should also consider the degree of absorption when determining the attractiveness between fireflies.

## C. Hybrid FOA-FA Algorithm

This section provides a comprehensive overview of the proposed algorithm, FOA-FA. The primary objective behind the development of FOA-FA is to address the limitations of the original FOA. The original FOA faces challenges in handling the negative domain, as it cannot generate candidate solutions uniformly across the problem domain. Moreover, it tends to prematurely converge due to the random term in Eq. (3), which typically produces small values within a radius of one around the best location. The methodology of the FOA-FA algorithm involves two distinct phases. The first phase makes use of FOA [22], where a group of fruit flies navigates in multiple directions using the ARM (Artificial Fruitfly Recognition Module). Consequently, these movements follow a uniform distribution across the problem space. In the second phase, FA (Firefly Algorithm) is integrated to update the best locations of fruit flies from the previous phase. This integration is essential to prevent FOA from getting stuck in premature convergence by combining its exploitation and exploration capabilities. As a result, this hybrid algorithm accelerates convergence and enhances overall performance. The primary steps of the proposed algorithm are outlined as follows:

Step 1: Initialization

- Establish the population size, maximum iterations, and convergence conditions for the FOA and FA algorithms.

- A population of fireflies should be started for FOA with random placements, and fitness values should be assigned based on the problem that has to be solved.

- Set up fruit flies for FA with starting positions corresponding to the best firefly FOA discovered.

- Decide on 0 iterations.

Step 2: Main Loop

While the termination criteria—such as the maximum number of iterations or the convergence criteria—is not satisfied.

Step 3: FOA Phase

- Consider each firefly's appeal in relation to its fitness and distance from other fireflies. Attractiveness increases with fitness level and distance travelled.

- Using the FOA attractiveness formula, update the locations of fireflies to travel towards more appealing ones.

- Reassess the changed roles' suitability.

- Based on fitness, choose the best firefly.

Step 4: FA Phase

- Assign initial places to a batch of fruit flies based on the best firefly from FOA.

- Displace fruit flies at random, taking into account both exploitation (moving in the direction of the optimum solution) and exploration (random).

- For every fruit fly, determine the fitness of the perturbed positions.

- Based on fitness, choose the finest fruit fly.

Step 5: Integration

- Compare the best fitness determined by FA and FOA.

- Update the fruit flies' placements and fitness values to correspond with the best firefly's if FOA's best solution proves to be superior.

- Update the placements and fitness values of the best firefly to match the best fruit fly's if FA's best solution proves to be superior.

Step 6: Termination

- Increase the number of iterations.

- Examine the criteria for termination. If it is satisfied, break out of the loop; if not, go back to step 2.

Step 7: Final Outcome

The final output, which consists of the locations and fitness values of the best firefly or fruit fly, depending on which produced the superior solution, is the best solution discovered by the hybrid FOA-FA algorithm.

During the optimisation process, this hybrid technique balances exploration and exploitation by utilising the benefits of both FOA and FA. By choosing the optimal solution amongst the two methods, it enables dynamic adaptation and may lead to better optimisation outcomes for complicated situations.

Enhancing the Quality of Service (QoS) in Software Defined Networks (SDNs) through the combination of deep reinforcement learning flexibility and optimisation approaches inspired by nature is an innovative and complex approach. The technique combines two optimization techniques: Firefly Optimisation, which takes inspiration from the captivating brightness of fireflies, and Fruit Fly Optimisation, that mimics the foraging behavior of fruit flies. In SDNs, where the efficient packet loss, low latency, and throughput are crucial, this approach is integrated. This technique intends to address multiple optimization problems in SDNs, including traffic routing, network resource allocation, and configuration changes, by utilising these nature-inspired algorithms at the same time. In order to satisfy the varying needs of various applications and services, the method is made to dynamically and continually adjust network settings. Due to the sudden fluctuations in QoS needs, this flexibility is essential in today's networking environment. Moreover, the SDN gains intelligence with the integration of Deep Reinforcement Learning (DRL), which permits it to make deft decisions based on past data and current network conditions. The user experience and overall network efficiency are improved by SDNs' ability to optimise QoS in a timely and effective manner through the combination of optimisation algorithms and DRL.

### D. Deep Q-Learning Framework

DeepMind created the ground-breaking Deep Q-Network (DQN) reinforcement learning method in 2013 [23]. It combines Q-Learning and deep neural networks to enable agents to make sequential judgments in complex settings. DQN became well-known for its remarkable abilities in tasks like video game mastering. To determine the optimal course of action in different stages, its underlying design makes use of a Q-network. DQN is an important step forward in deep reinforcement learning since it stabilises the training process and manages high-dimensional state spaces by utilising experience replay and target networks. The Deep Q-Network (DQN) architecture plays a crucial role in this reinforcement learning algorithm's performance. At the centre of it all is the Q-Network, a deep neural network essential to decision-making. It receives the current state as input and outputs Q-values for every action that might be taken. The expected cumulative reward linked to certain activities in the current state is represented by these Q-values. To implement the Q-

network, deep learning frameworks like as TensorFlow or PyTorch are frequently utilised. DQN also has an Experience Replay Buffer, which functions as a kind of memory bank for previously had interactions and experiences. Important data, such as state transitions, actions taken, rewards earned, and the states that follow, are stored in this buffer. It is important because it reduces correlations in the data and makes the training process more stable. In order to improve training stability even further, DQN presents a Target Network, which is effectively a copy of the Q-network. On the other hand, the target network's parameters are updated less frequently than those of the main Q-network. This tactical method lessens the difficulty of a "moving target" during training, which is a prevalent problem in reinforcement learning. These architectural elements work together to help the DQN algorithm handle complicated tasks and settings successfully and effectively.

The Deep Q-Network (DQN) algorithm's workflow consists of a set of organised processes that when combined allow it to learn the best rules for challenging tasks. Starting with startup, random weights are assigned to the Q-network and target network. Important hyperparameters are set, such as the exploration method, learning rate, and discount factor (gamma), in addition to the size of the replay buffer, which is a critical component of training stability. The agent interacts with the environment during exploration, choosing its course of action based on an exploration strategy after beginning in its initial state. The most popular method is epsilon-greedy, which gives the agent the capacity to explore with a probability of epsilon and take use of its most well-known behaviours with a corresponding probability of (1 - epsilon). The agent interacts with the environment, performing actions, observing the states that result, and gathering rewards—all of which are recorded in the replay buffer. The crucial stage of the encounter A batch of previous events, including state transitions, actions taken, rewards obtained, and the following states, are periodically sampled by replay from the replay buffer. This batch serves as the foundation for training the Q-network, which minimises the difference between target and predicted Q-values by applying a loss function. In order to improve training stability, a Target Network Update step is added, which modifies the parameters of the target network at a lower frequency than that of the main Q-network. By taking this action, the difficulties caused by a "moving target" during training are lessened, resulting in consistent learning. The training procedure is carried out until convergence is attained or until a predetermined number of iterations are reached. Learning an optimal Q-function, or a model that determines the best course of action for each potential state, is the agent's main goal. Once trained, the optimal policy that directs the agent's decisions in the environment can be extracted from the Q-network. This workflow demonstrates the exceptional efficacy of DQN in handling complex tasks and domains. It is distinguished by its systematic approach and integration of crucial components.

*1) Optimizing SDN QoS:* Firefly-Fruit Fly and DQN Fusion: The goal of improving Quality of Service (QoS) in Software-Defined Networks (SDN) has led to the creation of novel approaches, one of which combines Deep Q-Network

(DQN) with Firefly-Fruit Fly optimisation. This integration combines the best features of state-of-the-art deep reinforcement learning techniques with algorithms inspired by nature to optimise QoS and network performance in a fresh and promising way. Firefly-Vinegar Fly optimisation attempts to emulate the swarm intelligence of these natural organisms by taking cues from their collective behaviour. This allows the strategy to efficiently explore and make utilisation of network parameters, adapting and responding to shifting network conditions and user requests. By using a collective approach, the network may optimise routing and dynamically allocate resources, improving the quality of service overall. Deep Q-Network (DQN) is included to enhance the swarm intelligence and provide a higher level of sophistication to the decision-making process. Intelligent decision-making is made possible by DQN, which uses deep neural networks to assess and forecast the quality of potential actions in various network states. When these two approaches are combined, the network can quickly adjust to changing user demands, network conditions, and traffic volumes. This hybrid approach's potential to accelerate convergence and increase the network's responsiveness to real-time difficulties is one of its main advantages. By guaranteeing that resources are allocated optimally and lowering latency, packet loss, and other QoS-related problems, it also greatly increases the overall efficiency of SDNs. Essentially, there is a lot of room for improvement in QoS in SDNs because to the combination of DQN and Firefly-Fruit Fly optimisation technology. This approach, which makes use of deep learning and nature-inspired swarm intelligence, offers a flexible and adaptable way to deal with the difficulties and complexities of contemporary software-defined networks. This novel method opens the door to more effective and dependable communication as network needs increase and diversify, which makes it a potential direction for network optimisation and QoS improvement in the future.

## V. RESULT AND DISCUSSION

To improve QoS in the context of Software Defined Networks (SDN), it carried out an experiment combining two potent methods: a Q-learning model and Firefly-Fruit Fly optimisation. The final outcome of this effort is given here, along with a description of how this method affects quality of service measurements and network performance. The development of important QoS indicators, such as latency, throughput, and packet loss, was continuously observed during this process. Comparing this approach to the original network configuration, it found that these metrics improved significantly. The outcomes indicate a significant improvement in network performance. The emulator that was selected for the present research was MiniNet. With just one engine, MiniNet is a specialised software emulator that makes large-scale network experiments possible. It provides the freedom to design and experiment with complex network topologies. Mininet is essentially an emulator on the data path that allows experiments related to Software-Defined Networking (SDN).

### A. Latency in SDN Networks

Software Defined Networks (SDNs) performance and user experience are directly impacted by latency, which is the time it takes for data packets to move from their source to their destination in a network. Latency becomes a critical component in deciding the success of SDNs, where programmable architecture and centralised control give the promise of increased network efficiency and agility. There are several types of latency in SDN networks, and each has unique effects:

*1) Propagation latency:* The duration required for data packets to physically move across the network media is represented by this. Both the distance between the devices and the speed of light in the transmission medium have an impact on it.

*2) Transmission latency:* This is associated with the duration required to force data packets onto the medium of transmission. It is mostly reliant on the network devices' hardware and data rate.

*3) Processing latency:* The amount of time required for packet processing at the switches and controller greatly affects total latency in SDN settings. This entails tasks including rule matching, packet classification, and decision-making.

*4) Queuing Latency:* When packets build up in network device buffers, queuing happens. Packets that must wait in queue to be forwarded cause latency.

TABLE I. LATENCY IN SDN

| Network Load Level | Propagation Latency (ms) | Transmission Latency (ms) | Processing Latency (ms) | Queueing Latency (ms) |
|---|---|---|---|---|
| 2 | 1.2 | 0.5 | 0.1 | 0.2 |
| 4 | 1.5 | 1.0 | 0.3 | 0.4 |
| 6 | 2.0 | 2.5 | 0.5 | 1.2 |
| 8 | 2.5 | 3.2 | 0.8 | 1.8 |

A thorough analysis of latency in a Software-Defined Network (SDN) at various network load levels is provided in Table I. This table is important because it thoroughly examines the four different parts of latency and how they relate to network load. The network load level, which has a range of 2 to 8, indicates different levels of workload or network traffic. Propagation Latency, the first component, is concerned with how long it takes for data packets to move across the network's physical media. Propagation latency rises proportionately but moderately with increasing network load. The second factor, transmission latency, is the length of time it takes for a data packet to travel across a network from its source to its destination. As network demand increases, this component also suffers an increase.

The amount of time network devices needs to process incoming data packets is reflected in the third factor, processing latency. Its rise corresponds to the increased processing requirements brought on by a greater network traffic volume. Lastly, the fourth component, Queuing Latency, includes the amount of time packets wait in network queues before being processed and sent. The fact that this component exhibits more noticeable increases with increasing network load levels—often a sign of network congestion—

makes it especially interesting. Overall, Fig. 2 illustrates a refined understanding of the latency dynamics in SDNs. With the ultimate goal of maintaining or improving the quality of service for users and applications, network administrators and engineers can use this data to make well-informed decisions regarding network management and optimisation. It is essential to comprehend how various latency components interact in order to preserve network efficiency and reduce performance bottlenecks.



Fig. 2.    Latency in SDN networks.

### B. Jitter in SDN

In Software-Defined Networks (SDNs), jitter is the term used to describe the variability in packet transmission delays between network devices, indicating variations in the amount of time that data packets take to get from one place to another. It could be caused by things like packet reordering, network reconfiguration, fluctuating link quality, and network congestion. It can be problematic, especially for real-time applications. Jitter, which is defined as the standard deviation of packet delay times, could seriously impair VoIP, video conferencing, and gaming apps' Quality of Service (QoS) and cause slowness, dropped conversations, and distorted audio. Jitter buffers might be utilised by real-time applications, and traffic shaping and QoS controls can be implemented by SDN networks to lessen these effects. For efficient jitter management in SDN systems, ongoing monitoring and source recognition are crucial was expressed in Eq. (8).

$$Jitter = \frac{Overall\ Delay}{Overall\ Packets\ inbound - 1} \qquad (8)$$

TABLE II.    LATENCY IN SDN

| Network Load (Mb) | Jitter (ms) |
|---|---|
| 0.4 | 2.5 |
| 0.6 | 5.1 |
| 0.8 | 9.3 |

Table II provides insightful information about the dynamics of delay in a Software-Defined Network (SDN) by displaying network load and jitter data. Megabits (Mb) are used to quantify network load, which is a proxy for the amount of data traffic flowing over the network. Three different network load levels—0.4 Mb, 0.6 Mb, and 0.8 Mb—are shown in this table. Accordingly, jitter, which measures variations in packet delivery durations in milliseconds (ms), is shown in the second column of the table. Jitter is a crucial networking indicator. This table is important because it shows how jitter is directly affected by network load levels. Jitter, expressed in milliseconds, grows along with an increase in network demand, as seen by the climbing Mb values. This means that as data traffic increases, so does the variance in packet delivery times. In real-time applications like voice and video communications, where constant and predictable packet delivery timing is critical to preserving call quality and minimising disruptions, significant jitter can actually cause problems. Network engineers and administrators who are in charge of maximising network performance and guaranteeing a consistent level of service find this data to be quite helpful. Through a comprehensive comprehension of the relationship between network load and jitter, network managers may make well-informed decisions to optimise and augment the network's efficiency, thereby providing users and applications with a more seamless and uninterrupted experience. The jitter in SDN networks is graphically represented in Fig. 3.



Fig. 3.    Jitter in SDN networks.

### C. Throughput in SDN Network

Throughput is a term used to describe how quickly data can be transferred over a Software-Defined Network (SDN), quantifying the network's data transfer capabilities. Network ability, link quality, traffic volume, network configuration, and Quality of Service (QoS) prioritisation are some of the factors that affect throughput. Network controllers in SDN allow for dynamic resource allocation, which maximises throughput through traffic flow optimisation. While effective SDN management and scaling guarantee that applications and services receive the required data transfer rates for optimal performance, especially in bandwidth-intensive scenarios like video streaming and cloud services, accurate throughput

measurement is essential for assessing network performance was expressed in Eq. (9).

$$Throughput = \frac{Overall\ Packets\ Send}{Time\ of\ data\ send} \quad (9)$$

TABLE III. THROUGHPUT IN SDN

| Network Load (Mb) | Throughput (ms) |
|---|---|
| 5000 | 100 |
| 10000 | 50 |
| 15000 | 20 |

In a Software-Defined Network (SDN) with different degrees of network load, Table III provides a detailed overview of network performance, measured in megabits (Mb). Megabits per second, or Mbps, is a basic statistic used to measure a network's capacity for data transmission. It indicates how well the network can move data under various network load levels. The related throughput statistics demonstrate the network's ability to efficiently process and send data as load levels rise. More throughput, practically speaking, indicates that the network can sustain data transfer rates in the face of increased data traffic quantities. Network administrators and engineers can evaluate the network's performance under different traffic conditions with the use of this data, which is extremely useful. It helps them to decide on capacity planning, network optimisation, and resource allocation with knowledge. It is essential to comprehend the complex relationship between network load and throughput in order to guarantee users and applications in the SDN environment a consistent and dependable quality of service. The Throughput in SDN networks is graphically represented in Fig. 4.



Fig. 4.   Throughput in SDN networks.

### D. Packet Loss in SDN Network

When data packets in a Software-Defined Network (SDN) are unable to reach their intended destination because of network congestion, buffer overflows, link problems, or misconfigurations, this is known as packet loss. Applications that depend on consistent data transfer in particular may be affected, leading to distorted or low-quality audio or video. SDN networks use traffic engineering and Quality of Service (QoS) policies for path optimisation and prioritisation, and efficient buffer management lowers the chance of buffer

overflows. Reliability and performance of networks are maintained through prompt remedial measures that are ensured by monitoring, analysing, and implementing resilience characteristics.

TABLE IV. PACKET LOSS IN SDN

| Network Load Level | Packet Loss (%) |
|---|---|
| 40 | 0.5 |
| 80 | 2.0 |
| 120 | 5.0 |
| 160 | 10.0 |

Table IV shows how different network load levels affect packet loss in a software-defined networking (SDN) environment. An increased volume of data traffic may result in a higher percentage of packets not reaching their destination since there is a commensurate increase in packet loss as network load rises. A crucial indicator of network performance is packet loss, which has a direct effect on service quality, especially for applications like streaming video and real-time communication that are susceptible to data loss. Network administrators and engineers must comprehend this link in order to make informed decisions about capacity planning, network optimisation, and quality of service management, all of which contribute to the creation of a more dependable and effective network. The Packet Loss in SDN networks is graphically represented in Fig. 5.



Fig. 5.   Packet loss in SDN networks

### E. Accuracy Comparison

Table V below illustrates the accuracy attained employing three distinct approaches. Fig. 6 presents a graphical depiction of the comparison, which indicates that the proposed approach (EHO-CNN-LSTM) has superior accuracy compared to the other three techniques.

TABLE V. ACCURACY COMPARISON

| Technique | Accuracy |
|---|---|
| Navie Bayes | 88 |
| Neural Network | 81 |
| SVM | 78 |
| Proposed (DQN- FOA-FA) | 99.8 |

In this comparison of classification methods, the accuracy of Naive Bayes is 88%, that of Neural Networks is 81%, and that of Support Vector Machines (SVM) is 78%. With a remarkable accuracy of 99.8%, the suggested hybrid technique (DQN-FOA-FA) surpasses them all. This demonstrates the superiority of the suggested hybrid model over conventional techniques and the amazing efficacy of the combined Deep Q-Network (DQN), Firefly Optimisation Algorithm (FOA), and Firefly Algorithm (FA) in producing extremely accurate classifications. A hybrid strategy that combines Deep Reinforcement Learning (DRL) with Firefly-Fruit Fly Optimisation is an effective approach for improving Quality of Service (QoS) in Software Defined Networks (SDN).

*F. Discussion*

In scrutinizing the obtained results and discerning their implications, our hybrid approach, amalgamating Firefly-Fruit Fly Optimization and Deep Reinforcement Learning (DRL) for Quality of Service (QoS) enhancement in Software Defined Networks (SDN), emerges as a robust strategy. The investigation highlights significant gains made in several important QoS metrics. Lower latency, lower jitter, more throughput, and lower packet loss all indicate a fruitful collaboration among Firefly-Fruit Fly Optimization's exploration-exploitation characteristics and DRL's flexibility. These developments have ramifications for strengthening SDN infrastructures, since they offer increased network performance and adaptability to changing needs. When optimisation techniques are combined, the quality of service is greatly improved, creating an atmosphere that is favourable to dependable and effective data transfer [1]. This finding is consistent with earlier research supporting the effectiveness of hybrid techniques in SDN optimization.



Fig. 6. Accuracy comparison of suggested approach.

Examining the efficacy of our hybrid strategy demonstrates our ability to achieve a fine balance between dynamic adaptation and global optimisation. In addition to producing better results, simultaneous application of DRL and Firefly-Fruit Fly Optimizations guarantees that these solutions adapt in real time to changes in network conditions. When

compared to conventional optimisation methods, where the flexibility to accommodate dynamic network changes could be jeopardized, our approach's efficacy becomes evident [2]. Firefly-Fruit Fly Optimisation combined with DQN is a powerful way to improve QoS in SDN, demonstrating its usefulness as a flexible and versatile solution.

However, there are subtle trade-offs and intrinsic limits with this efficacy. Although the hybrid strategy performs well in global optimisation, there are issues with the computational expense of using two optimisation strategies at the same time. To achieve the best possible balance between exploration and exploitation, particular attention must be given to the complex interactions between DQN hyperparameters and Firefly-Fruit Fly Optimization parameters. Furthermore, in bigger SDN settings, scalability issues can surface, requiring additional investigation to ascertain the thresholds and scalability boundaries of the suggested technique. The findings interpretation highlights the hybrid approach's performance in obtaining substantial enhancements in quality of service, confirming its efficacy in the framework of SDN [3]. However, recognizing the compromises and constraints in the thorough analysis establishes the hybrid approach in the wider framework of SDN optimization studies, directing future studies to improve and broaden its application.

## VI. CONCLUSION

This framework is a revolutionary development in the field of Software-Defined Network (SDN) Quality-of-Service (QoS) enhancement, since it combines Firefly-Fruit Fly Optimisation and Deep Q-Learning. Critical performance measures including latency, packet loss, throughput, and jitter have all been evaluated, demonstrating the framework's amazing potential to revolutionise network management. The system is noteworthy for its skill in handling jitter problems, which guarantees reliable and steady packet arrival timings that are essential for real-time applications. Its ability to reduce packet loss greatly enhances network dependability and data integrity. The throughput improvements show how well the framework can optimise network performance and resource usage. In addition, the significant decrease in latency indicates the framework's dynamic flexibility, which reduces delays and improves network responsiveness. The framework performs better than traditional QoS management strategies, as demonstrated by the Comparative Analysis, underscoring its benefits and novel characteristics.

## VII. FUTURE WORK

Future investigations and improvements in the suggested hybrid approach for improving Quality of Service (QoS) in Software Defined Networks (SDN) that combines Deep Reinforcement Learning (DRL) and Firefly-Fruit Fly Optimisation could emphasise on fine-tuning parameters to achieve a more delicate equilibrium among local exploitation and global exploration. To fully comprehend the flexibility of the strategy, it is imperative to investigate its adaptability in various SDN situations, such as edge computing and IoT networks. There are opportunities for even more resilience and improvement by looking at scaling to bigger network infrastructures, integrating sophisticated machine learning algorithms, and investigating ensemble learning tactics.

Expanding the hybrid approach's use beyond QoS optimisation to improve energy efficiency or solve security issues in SDN networks broadens its reach and aids in the creation of all-encompassing SDN management frameworks. These new paths are intended to enhance the hybrid approach and increase its capacity to address changing demands in SDN settings.

## REFERENCES

[1] E. H. Bouzidi, A. Outtagarts, and R. Langar, "Deep reinforcement learning application for network latency management in software defined networks," in 2019 IEEE Global Communications Conference (GLOBECOM), IEEE, 2019, pp. 1–6.

[2] "Software Defined Networking: What is SDN?," Nutanix. Accessed: Nov. 17, 2023. [Online]. Available: https://www.nutanix.com/info/software-defined-networking

[3] R. S. Alonso, I. Sittón-Candanedo, R. Casado-Vara, J. Prieto, and J. M. Corchado, "Deep reinforcement learning for the management of software-defined networks and network function virtualization in an edge-IoT architecture," Sustainability, vol. 12, no. 14, p. 5706, 2020.

[4] M. Ye, J. Zhang, Z. Guo, and H. J. Chao, "Date: Disturbance-aware traffic engineering with reinforcement learning in software-defined networks," in 2021 IEEE/ACM 29th International Symposium on Quality of Service (IWQOS), IEEE, 2021, pp. 1–10.

[5] "Software Defined Networking (SDN): Benefits and Challenges of Network Virtualization - javatpoint." Accessed: Nov. 17, 2023. [Online]. Available: https://www.javatpoint.com/software-defined-networking-sdn-benefits-and-challenges-of-network-virtualization

[6] C. Yu, J. Lan, Z. Guo, and Y. Hu, "DROM: Optimizing the routing in software-defined networks with deep reinforcement learning," IEEE Access, vol. 6, pp. 64533–64539, 2018.

[7] Y. Li and Y. Qin, "Real-Time Cost Optimization Approach Based on Deep Reinforcement Learning in Software-Defined Security Middle Platform," Information, vol. 14, no. 4, p. 209, 2023.

[8] T. Yang, J. Li, H. Feng, N. Cheng, and W. Guan, "A novel transmission scheduling based on deep reinforcement learning in software-defined maritime communication networks," IEEE Transactions on Cognitive Communications and Networking, vol. 5, no. 4, pp. 1155–1166, 2019.

[9] A. Al-Jawad, I.-S. Comşa, P. Shah, O. Gemikonakli, and R. Trestian, "An innovative reinforcement learning-based framework for quality of service provisioning over multimedia-based sdn environments," IEEE Transactions on Broadcasting, vol. 67, no. 4, pp. 851–867, 2021.

[10] P. Sun, Z. Guo, J. Li, Y. Xu, J. Lan, and Y. Hu, "Enabling scalable routing in software-defined networks with deep reinforcement learning on critical nodes," IEEE/ACM Transactions on Networking, vol. 30, no. 2, pp. 629–640, 2021.

[11] P. Zhou et al., "QoE-aware 3D video streaming via deep reinforcement learning in software defined networking enabled mobile edge computing," IEEE Transactions on Network Science and Engineering, vol. 8, no. 1, pp. 419–433, 2020.

[12] I. Sarkar, M. Adhikari, S. Kumar, and V. G. Menon, "Deep reinforcement learning for intelligent service provisioning in software-defined industrial fog networks," IEEE Internet of Things Journal, vol. 9, no. 18, pp. 16953–16961, 2022.

[13] M. B. Hossain and J. Wei, "Reinforcement learning-driven QoS-aware intelligent routing for software-defined networks," in 2019 IEEE global conference on signal and information processing (GlobalSIP), IEEE, 2019, pp. 1–5.

[14] D. M. Casas-Velasco, O. M. C. Rendon, and N. L. da Fonseca, "Intelligent routing based on reinforcement learning for software-defined networking," IEEE Transactions on Network and Service Management, vol. 18, no. 1, pp. 870–881, 2020.

[15] E. H. Bouzidi, A. Outtagarts, R. Langar, and R. Boutaba, "Deep Q-Network and traffic prediction based routing optimization in software defined networks," Journal of Network and Computer Applications, vol. 192, p. 103181, 2021.

[16] D. M. Casas-Velasco, O. M. C. Rendon, and N. L. da Fonseca, "DRSIR: A deep reinforcement learning approach for routing in software-defined networking," IEEE Transactions on Network and Service Management, 2021.

[17] G. Kim, Y. Kim, and H. Lim, "Deep reinforcement learning-based routing on software-defined networks," IEEE Access, vol. 10, pp. 18121–18133, 2022.

[18] M. U. Younus, M. K. Khan, M. R. Anjum, S. Afridi, Z. A. Arain, and A. A. Jamali, "Optimizing the lifetime of software defined wireless sensor network via reinforcement learning," ieee access, vol. 9, pp. 259–272, 2020.

[19] H. P. Nugroho, M. Irfan, and A. Faruq, "Software Defined Networks: a Comparative Study and Quality of Services Evaluation," SJI, vol. 6, no. 2, pp. 181–192, Dec. 2019, doi: 10.15294/sji.v6i2.20585.

[20] H. Huang et al., "A new fruit fly optimization algorithm enhanced support vector machine for diagnosis of breast cancer based on high-level features," BMC Bioinformatics, vol. 20, no. 8, p. 290, Jun. 2019, doi: 10.1186/s12859-019-2771-z.

[21] "Firefly algorithm," Wikipedia. Aug. 08, 2023. Accessed: Oct. 24, 2023. [Online]. Available: https://en.wikipedia.org/w/index.php?title=Firefly_algorithm&oldid=1169297057

[22] X.-S. Yang, "Firefly Algorithms for Multimodal Optimization," in Stochastic Algorithms: Foundations and Applications, O. Watanabe and T. Zeugmann, Eds., in Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2009, pp. 169–178. doi: 10.1007/978-3-642-04944-6_14.

[23] J. Fan, Z. Wang, Y. Xie, and Z. Yang, "A Theoretical Analysis of Deep Q-Learning." arXiv, Feb. 23, 2020. Accessed: Nov. 17, 2023. [Online]. Available: http://arxiv.org/abs/1901.00137

# Revolutionizing Magnetic Resonance Imaging Image Reconstruction: A Unified Approach Integrating Deep Residual Networks and Generative Adversarial Networks

Dr M Nagalakshmi[1], Dr. M. Balamurugan[2], Dr. B. Hemantha Kumar[3], Lakshmana Phaneendra Maguluri[4], Dr. Abdul Rahman Mohammed ALAnsari[5], Prof. Ts. Dr. Yousef A.Baker El-Ebiary[6]

Associate Professor, Marri Laxman Reddy Institute of Technology and Management, Dundigal, Hyderabad-500043[1]
Associate Professor, Department of Computer Applications, Acharya Institute of Graduate Studies, Bengaluru[2]
Professor, Dept of Information Technology, RVR & JC College of Engineering, Guntur, AP, India[3]
Associate Professor, Dept. of Computer Science and Engineering,
Koneru Lakshmaiah Education Foundation, Vaddeswaram, Guntur Dist., Andhra Pradesh - 522302, India[4]
Department of Surgery Salmanyia Hospital Bahrain[5]
Faculty of Informatics and Computing, UniSZA University, Malaysia[6]

*Abstract*—Advancements in data capture techniques in the field of Magnetic Resonance Imaging (MRI) offer faster retrieval of critical medical imagery. Even with these advances, reconstruction techniques are generally slow and visually poor, making it difficult to include compression sensors. To address these issues, this work proposes a novel hybrid GAN-DRN architecture-based method for MRI reconstruction. This approach greatly improves texture, boundary characteristics, and picture fidelity over previous methods by combining Generative Adversarial Networks (GANs) with Deep Residual Networks (DRNs). One important innovation is the GAN's all-encompassing learning mechanism, which modifies the generator's behaviour to protect the network against corrupted input. In addition, the discriminator assesses forecast validity thoroughly at the same time. With this special technique, intrinsic features in the original photo are skillfully extracted and managed, producing excellent results that adhere to predetermined quality criteria. The Hybrid GAN-DRN technique's effectiveness is demonstrated by experimental findings, which use Python software to achieve an astounding 0.99 SSIM (Structural Similarities Index) and an amazing 50.3 peak signal-to-noise ratio. This achievement is a significant advancement in MRI reconstructions and has the potential to completely transform the medical imaging industry. In the future, efforts will be directed towards improving real-time MRI reconstruction, going multi-modal MRI fusion, confirming clinical effectiveness via trials, and investigating robustness, intuitive interfaces, transferable learning, and explanatory techniques to improve clinical interpretive practices and adoption.

*Keywords—Magnetic Resonance Imaging (MRI); deep learning; generative adversarial network; deep residual network; ResNet50*

## I. INTRODUCTION

Body organs, extremities, and different tissues were imaged using medical imaging devices such Magnetic Resonance Imaging (MRI), ultrasound, computed tomography (CT), and X-ray. Nevertheless, poor signal-to-noise ratios (SNR) and reduced contrast-to-noise ratios (CNR), as well as image artefacts, may be present in images obtained using such imaging paradigms. Especially in clinical circumstances, the impact of imperfections and noise from many sources, like those caused by magnetic fields inhomogeneities, have to be matched with the period needed for collection. Actually, since CT scans have been so much speedier than MRI scans, they are frequently chosen over MRI notwithstanding its weaker soft tissue contradiction, loyalty, and the presence of ionising energy. The average MRI session lasts between 20 and an hour for each sufferer. Contrasted to CT, individual scans take longer, and they frequently call for distinct MR patterns that react distinctly to the properties of distinct tissue kinds such as T2-weighted (T2-w) or T1-w MRI [1]. To solve these difficulties and enhance the images quality for improved visual perception, comprehension, and assessment, image reconstruction methods were created. Radionics, medical image assessment, and computer-assisted identification and assessment are a few applications of deep learning (DL) techniques that were employed effectively in medical imaging [2], [3]. During the past ten years, enthusiasm for artificial intelligence (AI) had significantly increased across almost all branches of research and innovation. The development of DL, a collection of techniques centred around artificial neural networks (ANN), which has shown to be a powerful all-purpose tool for automated processes, has played a significant role in its rise to prominence [4].

As non-linear models of the information that best match the intended objective for which it was learned, the DL networks learn the fundamental features and significant basis functionalities. It performs this completely autonomously (i.e., without operator input) as a component of an improvement procedure to identify the best attributes. In contrast to conventional statistical and machine learning (ML)

techniques, DL algorithms performance improves accordingly to a power law as additional information is added [5]. Convolutional neural networks (CNNs), one of the DL approaches, have transformed imaging and computer vision since they offer a data-driven strategy for addressing a variety of difficult issues. Due to their data-driven type, DL techniques frequently outperform conventional linear analytic techniques, where these characteristics have been manually chosen and created over extensive experience and research. The term "deep" typically refers to the dimensions and quantity of "layers" that make up a neural network, every one of which has the capacity to store an enormous amount of these properties. Acquiring the interconnected neuron's weights from the initial strata to all the following ones until the final strata has been a necessary step within the optimisation procedure [6]. This procedure is typically labelled or supervised in advance. By making modest adjustments to the results, or "returning" of the system, stochastically gradient descent has been utilised during learning to determine this network weight of neurons that are optimum. Backpropagation constitutes a technique whereby inefficiencies from the identified and anticipated outcome can be reduced by changing the weights [7]. As a result, the artificial neurons can successfully generate non-linear judgements, and the system design can parameterize the problem and automatically choose features by controlling the weights among the neurons. Among the key factors driving the AI development researchers are seeing today was the development of novel networks, network structures, and neuronal connections [8].

Recently, DL, also known as representation learning [9], had received a lot of interest for the evaluation of medical images [10]. Deep learning (DL) outperformed traditional machine learning (ML) approaches because it can extract features from unfiltered data sources during the learning phase. It may acquire concepts based on inputs owing to its multiple secret layers [11]. The DL technique and its usage in several fields [2], Reconstruction of medical images has been aided by recent improvements in successful computing platforms such as cloud-based computing and graphics processing units (GPUs). The quantity of research carried out on medical image recovery has expanded dramatically in recent years. Acceleration of MRI scans is one of these applications. This problem is essential because, while MRI remains the leading diagnostic technology for a range of procedures, it is inherently slower than other methods due to the physical processes involved in data acquisition. Thus, a critical component in the MRI widespread clinical use was the lowering of scan times [12].

Conventional image reconstruction techniques have traditionally relied on broad (unsupervised) restoration techniques that make very few implications about the object being photographed. These approaches do not conveniently support methods that would greatly boost up MR deals, despite the fact that they give users trust in the pictures by providing solutions that have been typically resilient to noise and distortions. The speed factors that can be achieved for anatomical contrasting images, like T2-w and T1-w images, without taking into account prior understanding have been in the range of 2 to 4. The quality of the restoration then quickly deteriorates with obvious artefacts. Greater acceleration variables may be attained for operational contrast, like coronary MR angiography (CMRA) [13].

New prospects for dramatically increasing MR image capture and restoration speeds whilst retaining high quality have emerged with the introduction of DL-based techniques into MRI restoration. Sparse restoration, multi-contrast restoration, and parallel imaging are the three key areas where AI-based MR imaging has made progress. Compressed sensing (CS), a method for quick MR imaging predicated on the sparseness earlier, has gained prominence in recent years. Nevertheless, the repetitive solution process requires a considerable amount of effort to produce a high-quality restoration, and the regularisation parameter has been chosen empirically. Despite the fact that several numerical techniques, including Stein's Unbiased Risk Estimation (SURE) [14], were suggested to optimise the free variables in MRI, these techniques are saddled with significant computational expenses.

Additionally, just a small number of standard images are used in most systems, and the preceding information is.

The primary traditional challenge in MRI image reconstruction is to accurately reconstruct high-quality images from raw data acquired during the scanning process, while addressing issues such as under sampling artifacts, noise, and motion-related distortions [15]. These factors can lead to blurry, distorted, or incomplete images, hindering accurate diagnosis and interpretation by medical professionals. Overcoming these challenges is crucial for enhancing the diagnostic value and overall effectiveness of MRI imaging in clinical settings. The significance of this study lies in its innovative approach to addressing the challenges of MRI image reconstruction, which is crucial for accurate medical diagnosis and monitoring. While MRI is a valuable diagnostic tool, its relatively slow imaging speed can limit its practicality. This paper introduces a Hybrid Deep Learning framework that combines the strengths of Generative Adversarial Networks (GANs) and ResNet50 Deep Residual Networks to enhance the quality and speed of MRI image reconstruction. By effectively leveraging GANs for image refinement and ResNet50 for feature extraction, the proposed approach not only accelerates the reconstruction process but also ensures improved image fidelity, texture, and edge preservation. This has the potential to greatly benefit clinical practice by enabling quicker and more accurate interpretations of MRI scans, thus enhancing the overall efficiency and reliability of medical imaging for diagnosing and managing various illnesses. The study's utilization of performance metrics such as PSNR and SSIM provides quantitative evidence of the effectiveness of the proposed method, further underscoring its importance in advancing MRI image reconstruction techniques beyond the limitations of previous approaches. Through extensive experimentation and comparisons with existing methods, the research demonstrates substantial improvements in MRI image reconstruction. In prior MRI image reconstruction research, a notable issue was the inherent trade-off between image quality and reconstruction speed. Traditional methods struggled to provide

high-quality images in real-time or with limited computational resources. This issue posed challenges in clinical applications where timely and accurate diagnoses are imperative. The study "Improvement of MRI Image Reconstructions through Combination of Deep Residual Networks and Generating Adversarial Networks in a hybrid Deep Learning Architecture" made the following significant contributions:

- Better Image Quality: By incorporating GANs, MRI images may be seen more clearly and have greater value as diagnostics.

- Real-time or Near-real-time Reconstructions: Applications requiring MRI reconstruction of images in actual time or near-real-time are critical in medical applications. Deep ResNets facilitate rapid MRI reconstruction of images.

- Balancing Trade-off: This study overcomes the enduring difficulty of choosing one as opposed to another in earlier approaches by striking an equilibrium between the quality of images as well as reconstruction speed.

- Clinical Utilisation: By rapidly delivering images of excellent quality, the suggested hybrid deep learning system will greatly enhance medical evaluations and support healthcare providers in making quick and precise decisions.

- Improvement of MRI Technologies: This study advances MRI technology by utilising the advantages of GANs as well as ResNets, guaranteeing that patients obtain the greatest imaging results and medical experiences.

The rest of the paper is structured as follows: Section II gives an overview of relevant studies. Section III covers the approach's research gap. The proposed approach is presented in Section IV, which describes the improving MRI Image Reconstruction by employing a Hybrid Deep Learning Mechanism based on Generative Adversarial Networks with Deep Residual Networks. Section V goes over the findings and performance analysis. Finally, Section VI summarises the contributions of research techniques to the work

## II. RELATED WORKS

It is difficult to recreate a PET image regarding low-count projected data and physical consequences since the inverse issue is poorly stated and the final image has been typically noisy. GANs were recently demonstrated improved performance in a variety of computer vision applications, which has sparked significant attention in medical diagnostics. For the purpose of reducing streaking aberrations and enhancing the PET images quality, Qianqian et al. [15] suggested a unique deep residual GAN (DRGAN) framework relying on GANs. Instead of directly generating PET images, the investigators of the suggested approach taught a generator to build a "residual PET map" (RPM) for expressing images. For the purpose of requiring anatomically accurate RPM and PET images, DRGAN employed two barriers (critics). The authors created residual dense links with pixel shuffle processes (RDPS blocks), which promote feature reusing and

avoid resolution loss, to better enhance the contextual data. To assess the suggested strategy, both simulated data and actual medical PET data have been employed. The quantification outcomes demonstrate that DRGAN may achieve higher performances in the bias-difference trade-off and deliver equivalent image quality when compared to other cutting-edge algorithms. The comprehensibility of the produced residual PET map (RPM) or the topographical precision of the restored PET images are not included in the description. For accurate medical evaluation and clinical choice-making, anatomy precision is just as important as comprehension. Evaluation of the quality and clinical applicability of the restored PET scans is crucial.

Using sparsely collected information to speed up MRI causes substantial artefacts that make it difficult to see the image's actual content. Fully dense attention CNN (FDA-CNN), a cutting-edge convolutional neural network (CNN), was suggested by Biddut et al. [16] to eliminate aliasing artefacts. Incorporating fully dense connection and an attention procedure for MRI restoration, researchers upgraded the Unet framework. The fundamental advantage of FDA-CNN was that every decoder strata's attention gate boosts learning by concentrating on the pertinent picture information and improves network generalisation by eliminating irrelevant engagement. Convolutional layers with close connections can also reuse map features and avoid the disappearing gradient issue. The investigator additionally employs a fresh, effective under-sampling sequence in the phase orientation that extracts high and low frequency bands both arbitrarily and deliberately from the k-space. Three distinct datasets and sub-sampling masks were used to objectively and subjectively assess FDA-CNN's performance. The suggested method outperformed five existing DL-oriented and 2-compressed sensing MRI restoration methods because it produced cleaner and more illuminated images. Additionally, FDA-CNN outperformed Unet for an accelerating value of 5 in terms of average SSIM, VIFP, and PSNR each by 0.35, 0.37, and 2 dB, respectively. The calibre and variety of the training databases have a significant impact on the FDA-CNN technique's success. The network's capacity to generalise may be hampered if the training databases have been too small or do not sufficiently cover all of the potential MRI changes.

A unique framework was created by Manimala et al. [17] to quickly and accurately restore chaotic sparse k-space information into MRI. In the suggested approach, Rician noise-corrupted MR pictures are denoised using a CNN. The technique takes advantage of signal similarities by analysing related patch collectively in order to extract local information. The CNN was learned on a GPU with the Convolutional Network for Fast-Feature Embedding system, rendering it appropriate for online restoration, resulting in a crucial drop in running time. The primary benefit of the CNN-based restoration over current cutting-edge methods has been the elimination of the need for noise level optimisation and forecasting during denoising. Different under-sampling strategies were used in analytical studies, and the findings show great accuracy and a constant peak SNR especially at twenty-fold under-sampling. Large under-sampling rates make it possible to transmit k-space data wirelessly, and fast

restoration makes it possible to use our approach for virtual health surveillance. The fact that the CNN had been trained on a GPU according to the research suggests that there may be significant processing demands. In pragmatic medical situations, it is crucial to take into account the system's scalability and hardware capabilities required for real-time or online reconstruction.

The amount of time it takes to acquire data can be significantly decreased with CS-MRI. The conventional CS-MRI approach, which relies on iterations, is versatile in modelling but typically time-consuming. Because of its excellent effectiveness, the deep neural network (DNN) approach has recently gained popularity in CS-MRI. The DL method's disadvantage, however, is its rigidity. It heavily relies on the k-space data's screening process and learning images. In order to achieve speedy, adaptable, and accurate restoration, Ruizhi and Fang [18] suggested an iterative technique for MRI restoration termed IDPCNN that combines the advantages of both the conventional approach and the DL techniques. Projection and denoising are two stages of the suggested approach. A cutting-edge denoiser has been used in the denoising process for smoothing the images. The projection stage continuously adds specifics to the space domain while exploring the previously acquired frequency domain data. Under various sampling filters and rates, the restoration quality has been better than the finest MRI restoration techniques. The IDPCNN offers the possibility for broad clinical uses due to its stability, rapidity, and high restoration quality. The study makes no reference to any clinical verification or assessment of the IDPCNN technique by qualified physicians or medical professionals. To ascertain whether the reconstructed images have been diagnostically precise and trustworthy for generating medical judgements, clinical evaluation is essential.

Dynamic MRI is a valuable technique for capturing the ever-changing anatomy of various organs in the human body over time. However, its clinical utility is often constrained by practical limitations, such as limited acquisition time due to mechanical and physiological factors. Dynamic MRI is known to exhibit spatio-temporal heterogeneity in its frequency spectrum, particularly in the k-space domain. To address this challenge and expedite the acquisition process, researchers Shashidhar and Subha [19] devised a novel approach involving a cascaded Convolutional Long Short-Term Memory (ConvLSTM) framework. This technique focuses on restoring T2-weighted dynamic MRI patterns from significantly under-sampled k-space data. Specifically, it leverages a Cartesian undersampling mask to acquire less k-space data than traditionally required. The ConvLSTM framework plays a pivotal role in mitigating aliasing artifacts resulting from this undersampling process. Notably, it excels in capturing both temporal and spatial connections within the image data, surpassing the capabilities of conventional CNN-based restoration methods. While the use of medical databases, such as the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset, offers valuable insights, it also raises ethical concerns related to data privacy and informed consent. It is imperative that research endeavors adhere to rigorous ethical standards and secure the necessary permissions and authorizations for the responsible use of such data.

In the areas of neuroscience and ML sectors, there has been an increase in interest in comprehending how the individual brain functions. In earlier research, generative adversarial networks (GANs) and autoencoders were used to boost the accuracy of stimulus image restoration from functional MRI (fMRI) datasets. These approaches, however, primarily concentrate on gathering pertinent aspects across two separate data modalities, namely, fMRI and stimulus images, whereas neglecting the fMRI statistic's temporal information, which results in less than the ideal efficiency. Shuo et al. [20] suggested a temporal information-guided GAN (TIGAN) to restore visual data from the brain's activity to tackle this problem. The suggested approach is made up of three essential parts, particularly: an algorithm for image restoration that is employed to render the restored image more comparable to the unique image; a fMRI encoder for visualising the stimulus visuals into hidden space; along with an LSTM framework for fMRI characteristic mapping. Additionally, researchers use a pairwise ordering loss to organise the stimulus images as well as fMRI to guarantee strongly correlated pairings are at the highest level and poorly associated ones are at the lower level in order to better quantify the relationship between two distinct types of databases (i.e., natural images and fMRI). The empirical findings on real-time databases imply that the recommended TIGAN outperforms a number of cutting-edge image restoration techniques. GANs and LSTM architectures may be operationally taxing; learning and inference need a significant amount of computational energy and time. The TIGAN method's computational effectiveness is not mentioned in the outline, which can have an impact on its accessibility and flexibility, especially in large-scale or real-world applications.

A DL-based restoration approach was presented by Yan Wu et al. [21] to enhance image quality for fast MRI. Researchers combined a volumetric recursive deep residual CNN with the self-attention process, which recorded long-range connections across image areas. Each convolutional layer included a self-attention component integrated into it, which computed the signal at each point as the weighted average of all the attributes at each position. Additionally, data consistency has been mandated, and reasonably dense shortcut links have been used. The SAT-Net suggested network has been implemented to cartilage MRI data that was obtained employing an ultrashort TE pattern and retroactively under-sampled in a pseudo-random Cartesian sequence. The algorithm had been evaluated employing 24 images that produced better results after being learned on 336 3-dimensional images (each one comprising 32 slices). The structure is adaptable to a wide range of applications. Although it is said that the SAT-Net approach can be applied to a variety of applications, the explanation does not give any concrete examples or evidence of its usefulness outside of cartilage MRI. Validating the approach's efficacy and adaptability with diverse MRI data kinds and clinical uses is crucial.

To speed up parallel MRI with residual complicated CNN, Shanshan Wang et al. [22] introduced a multi-channel image

restoration technique called DeepcomplexMRI. To learn the deep residual CNN offline, DeepcomplexMRI employs an extensive amount of previous multi-channel ground truth images as emphasise data, in contrast to most current efforts that depend on the utilisation of coil sensitivity or prior knowledge of predetermined transforms. A sophisticated convolutional network has been specifically suggested to account for the relationship between actual and fictitious portions of MRI. Additionally, among network levels, the k-space information coherence is continually guaranteed. The suggested technique can recover the intended multi-channel pictures, according to tests using in vivo databases. Its contrast with cutting-edge techniques also shows that the suggested method could more precisely rebuild the intended MR images. The utilisation of medical imaging database poses ethical questions about data usage, privacy, and informed permission. It is crucial to confirm that the research complies with ethical standards and has gotten the necessary authorization for the gathering and use of data.

The main problem with the aforementioned techniques is that they require more thorough assessment and validation with respect to both clinical application and technical excellence. These papers offer novel deep learning-based solutions for a range of healthcare imaging problems, such as MRI and PET imaging restorations; but they frequently do not have comprehensive medical confirmation, expert medical confirmation, and scalable evaluations. The evaluated methods suggest novel approaches to improving medical images. Though DRGAN emphasizes PET image improvement with residual GANs, questions arise over the clinical relevance due to inadequate focus on RPM comprehension and anatomical correctness. FDA-CNN covers MRI artifact reduction, although problems with scalability and real-time equipment needs are not explored. CNN-based restorations for unpredictable sparse MRI data offers potential for denoising but require investigation into scalability and real-world applications. IDPCNN provides excellent CS-MRI restoration, although the lack of clinical validation raises questions regarding diagnostic accuracy. ConvLSTM for continuous MRI raises ethical issues without providing details on adherence to norms. TIGAN provides fMRI recovery with temporal instructions, although the computational efficacy is not explored. SAT-Net produces promising outcomes for rapid MRI enhancement, but its wider relevance remains unsubstantiated. Deepcomplex MRI presents a unique approach for simultaneous MRI reconstructions with recognized ethical implications, emphasizing the importance of validation over several datasets and clinical contexts. Moreover, informed permission and information protection are two moral problems that have not been often handled [19]. The primary area of study deficiency is the inadequate use of these intriguing deep-learning approaches in realistic clinical contexts, where their efficacy, simplicity, and robustness must be thoroughly examined and verified. Furthermore, evaluating the actual worth and benefits of these suggested approaches in clinical settings is hampered by the dearth of studies compared with current state-of-the-art procedures.

## III. PROBLEM STATEMENT

Previously, CNN-based techniques had trouble maintaining subtle image information and materials, LSTM-based techniques had trouble capturing long-range dependence, and conventional machine learning techniques had trouble adapting to intricate image-to-image conversions in MRI reconstructions. To solve these problems, the suggested study proposes a combination of deep learning architecture that blends Deep ResNets with GANs. Through the use of GANs, this method improves texture retention and quality of images while facilitating real-time or almost real-time reconstruction—a crucial feature for applications in healthcare. This work expands the application of MRI reconstruction to medical imaging by crossing the speed and quality of image gaps, providing a more flexible and effective means of enhancing healthcare outcomes and increasing diagnostic precision.

## IV. PROPOSED GAN-RESNET50 APPROACH

Undersampling is a key method for speeding the capture of Magnetic Resonance (MR) images in Magnetic Resonance Imaging (MRI). It may shorten the duration and expense of MR scanning, minimizing patient pain. Generative adversarial network (GAN) [23] is a particularly effective MRI techniques for obtaining the original picture from undersampled images. In this work, an MRI picture and the accompanying ground truth images are used. Divide the data into training and testing sets after normalizing the intensity of pixels as part of the preprocessing. Create the ResNet50 generator and discriminator network-based GAN framework. Distinguishing between genuine high-quality MRI images and artificially produced high-quality MRI images is the goal of the discriminator. The goal of the generator is to provide excellent MRI images that can deceive the discriminator. Construct the GAN framework's loss function. A content loss must be included in the loss function together with adversarial loss, which pushes the generator toward generating realistic images. This guarantees that the generated images and the real-world images are similar. Iteratively improving the loss function will train the GAN model. The discriminator becomes more adept at distinguishing actual images from created ones, while the generator gets better at producing high-quality MRI images. Utilize quantitative criteria to assess the reconstructed images, such as peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM). These parameters assess how well the reconstructed pictures match the originals. The experimental results reveal that the hybrid GAN-DRN technique achieves amazing performance with an SSIM and a PSNR, resulting in a significant improvement in picture fidelity. These findings point to the possibility of transformational uses in MRI reconstructions. However, the technique's performance needs to be evaluated in a variety of clinical circumstances to confirm its durability and practical efficacy. Fig. 1 shows the overall methodology of the Proposed Method.

Fig. 1. Overall workflow of the proposed approach.

## A. Dataset

The study utilized the dataset from the MICCAI 2013 Grand Challenge to assess GAN-ResNet50 performance. T1 weighted coronal brain sections make up the dataset. For training, 6277 photos are used for validation, while 9901 images are used for testing. It presents the results of the test set in this work unless specified otherwise. Each image has a size of 256 by 256 pixels, with pixel values ranging from 0 to 1. In the experiment, the center of the original photos is cropped if the needed size is less than 256*256 in order to assess the model's ability to reconstruct images of various sizes, including 64*128, 128*256, and 256*256. To ensure that the number of photos needed for training, validation, and testing is the same for images of all sizes, thus only crop one patch per image is made. It implies that the results on photos of various sizes may be compared fairly. Additionally, when a model is tested, photos of the same size as those used to train it are used [24].

## B. Data Pre-processing

For categorization, the normalization type of pre-processing is essential. The input data should be normalized to speed up the learning process. Furthermore, to avoid numerical problems like accuracy loss due to arithmetic errors, some sort of data normalization may be necessary. Following initially outweighing features with originally lower ranges, characteristics with initially big ranges would take over a gradient descent. Feature space normalization might be considered a kernel impression of pre-processing rather than, strictly speaking, a type of pre-processing because it is not introduced externally to the input matrices. In other words, by transforming the data into a usable plane, normalization is a distinct kernel mapping approach that simplifies calculations. Given the enormous amount of data points, the complex normalization algorithm requires an extended period for processing. The Min-Max normalization technique that was selected is fast and efficient.

By using Min-Max Normalisation, the real data m is translated straight into the required interval as given in (1) $(max_{new}, min_{new})$.

$$m = min_{new} + (max_{new} - min_{new}) * \left(\frac{m - min_x}{max_x - min_x}\right) \quad (1)$$

The method's advantage is that it maintains every connection among the data points precisely.

## C. Image Segmentation

Image segmentation is a critical component of medical image processing, including MRI data enhancement. In the context of MRI reconstruction, image segmentation plays a pivotal role in isolating and delineating specific anatomical structures or regions of interest within the reconstructed

images. This segmentation process is fundamental for tasks such as tumor detection, organ volumetry, and the visualization of specific tissues or pathologies. The integration of GANs and Deep ResNets within MRI reconstruction pipelines enhances the segmentation process by providing high-quality, noise-reduced images for subsequent analysis.

### D. Preprocessing for Segmentation

Before the actual segmentation process, the reconstructed MRI images are subjected to preprocessing steps, which may include denoising, bias field correction, and intensity normalization. The integration of GANs and ResNets at earlier stages of image reconstruction inherently improves the quality of MRI data, leading to more accurate and robust preprocessing results.

Semantic and Instance Segmentation

- In the context of medical imaging, there are two primary types of image segmentation: semantic segmentation and instance segmentation.

- Semantic segmentation involves classifying each pixel in the image into predefined categories or labels. For example, it can be used to differentiate between different types of tissues, such as white matter, gray matter, and cerebrospinal fluid in brain MRI.

- Instance segmentation takes the process further by distinguishing individual instances of objects within the same category. For example, it can be used to identify and differentiate multiple tumors in an MRI scan.

*1) ROI Detection:* Region of interest (ROI) detection is a vital aspect of medical image analysis. GANs and ResNets contribute to the creation of sharper, high-resolution images, making it easier for automated algorithms to detect and segment ROIs accurately. This is especially valuable in tasks such as identifying specific pathologies or measuring anatomical structures.

*2) Benefits of GANs and ResNets in Segmentation:*

- Noise Reduction: The noise reduction capabilities of GANs help in segmenting images with minimal interference from artifacts or random noise.

- Edge Preservation: ResNets, with their deep architectures, are well-suited for preserving image edges, enabling the precise delineation of structures in the segmentation process.

- Better Contrast: GANs enhance the contrast and clarity of MRI images, making it easier for segmentation algorithms to distinguish between different regions or structures.

- Improved Generalization: By improving the overall quality of MRI images, the integration of GANs and

ResNets ensures that segmentation models generalize better across different datasets and clinical scenarios.

### E. Generative Adversarial Network with ResNet50

A generator $G_r$ and a discriminator $D_r$ make up a GAN. The parameters or weights of the generator are indicated as $G_r$ and $D_r$, respectively. The generator reconstructs the original, completely sampled MRI (I) and labels them as real, while the discriminator is instructed to identify the reconstruction from the undersampled images (n) as false, or not real. The discriminator cannot determine if the images it reconstructs are real or artificial, that is, it is unable to differentiate the reconstructed images apart from the completely sampled ones. In contrast, the generator is trained to attain the opposite goals. The following loss function is able to be used to mathematically represent the entire process of training as a minimax activity given in (2):

$$\mathcal{L}\left(\theta_{G_r}, \theta_{D_r}\right) = min_{\theta_{G_r}} max_{\theta_{D_r}} \mathbb{E}_{i \sim d(i)}[\log D_{r_\theta}(i)] \mathbb{E}_{n \sim d_n(n)}[\log(1 - D_{r_\theta}(G_{r_\theta}(i)))] \qquad (2)$$

Here the distribution of the undersampled images is represented by $d_n$ and the distribution of the fully sampled images by $d$. It has been demonstrated that the generator minimizes the Jensen-Shannon divergence among the distributions of the original and reconstructed images when using an optimized discriminator. In other words, it is possible to think about GAN models as minimizing the distributional difference between completely sampled and reconstructed images.

*1) Generator architecture with ResNet50:* The study utilized a unique generator design to enhance training stability and speed up model convergence. ResNet50 can capture scale-invariant features based on spectral data, convolutional processes are good at capturing spatial properties. Therefore, it is preferred to take into account both the spatial and spectral data in one framework. Their main distinction is that ResNet50 operates on each subband from the preceding level. The ResNet50 feature, according to the study, may speed up feature learning. As a result, it might provide quicker convergence and more training consistency. This is the first time a deep learning model has been tried using ResNet50-based architectures.

When using ResNet-50 as the generator in a GAN framework in Fig. 2 for MRI image reconstruction, the generator network plays a crucial role in transforming low-quality MRI images into high-quality ones. The ResNet-50 architecture consists of convolutional layers, residual blocks, downsampling layers, and an output classification layer. To repurpose ResNet-50 as a generator, the architecture is modified by removing the classification layer and adapting the last layer to match the dimensions of high-quality MRI images. The modified ResNet-50 generator takes low-quality MRI images (LQ) as input and aims to generate corresponding high-quality MRI images (HQ).

Fig. 2.    Generator architecture with ResNet50.

The input LQ images pass through the initial convolutional layers, capturing low-level image features. The residual blocks, which are the hallmark of ResNet architectures, help the generator learn the residual mapping between LQ and HQ images. These blocks contain skip connections that directly connect earlier layers with deeper layers, allowing the network to learn the difference or residual between the input and target images. By incorporating skip connections and residual learning, the generator can effectively capture and preserve important image details during the upsampling process. The upsampling layers gradually increase the spatial resolution of the input, helping to generate high-quality images with fine details. The modified last layer of the generator outputs high-quality MRI images that closely resemble the ground truth images. During the training process, the generator is updated based on the adversarial loss and content loss, as mentioned earlier. This encourages the generator to produce high-quality images that can deceive the discriminator and resemble the ground truth high-quality MRI images. By leveraging the ResNet-50 architecture as the basis for the generator, the GAN framework benefits from its ability to learn complex image mappings and effectively handle attribute extraction, resulting in improved MRI image reconstruction.

*2) Discriminator with ResNEt-50:* The flattening layer of the discriminator employed in GAN-ResNet50 is the primary factor enforcing the input size constraint. Considering the GAN-ResNet50 discriminator is analysing a 256×256-pixel MR picture. The image is compressed by 26 = 64 times, or 4×4 pixels, after undergoing 6 convolutional procedures. The 'concat, 1024' layer has 1024 channels when it is reached. Its size at this layer could be characterized as a tensor of 4×4×1024 dimensions. The tensor from the previous layer,

"concat, 1024," is unwound into a column vector of 8192 units by the following flattening layer. The following dense layer, which has 8192 input nodes, is capable of accepting it. On the other hand, if the similar model is given an image with 128*128 pixels. This image would be converted into a 2×2×1024-dimension tensor, unraveled into a 4096-unit column vector, and then layered with an 8192-unit dense layer even though there aren't enough input nodes for it. This instance demonstrates how the GAN-ResNet50 discriminator's input size limitation works.

To eliminate this constraint, the investigation replaces the flattened procedure with a global average pooling (GAP) layer. 29 GAP determines the average value of each pixel in each input way, irrespective of the image's size or pixel count. For instance, the GAP layer inputs in the simulation have 1024 channels. The result of this function is always going to be a series of vectors with 1024 units, no matter how big the tensor that was input from 'concat, 1024' is. Following that, the sigmoid is activated by applying an intense layer of 1024 units. As a result, the input size restriction is significantly lifted by the discriminator's architecture. According to the paper, this is a plan to use GAN-based MRI techniques for the first time.

While ResNet-50 is typically used for classification tasks, discriminator architecture in Fig. 3 can be designed to effectively assess the authenticity of the generated images. The discriminator architecture can be constructed by modifying the last layer of ResNet-50 to output a binary classification result, indicating whether the input image is real (high-quality) or fake (generated). By removing the classification layer and retaining only the convolutional layers and downsampled features of ResNet-50, the discriminator focuses on learning discriminative features for distinguishing between real and generated images. The modified ResNet-50 discriminator can have additional layers, such as fully connected layers and a final sigmoid activation layer for binary classification. The convolutional layers within the ResNet-50 architecture capture hierarchical image features, while the additional layers facilitate the final decision-making process. The discriminator takes input images, both real high-quality MRI images (HR) and the generated high-quality images (G(LQ)), and outputs the probability of the input being a real image. Throughout the training process, the discriminator undergoes training to minimize the binary cross-entropy loss, which measures the disparity between its predictions and the actual ground truth labels. This adversarial training strategy compels the discriminator to become proficient in discerning between authentic and synthesized images. Concurrently, the generator, which is based on the modified ResNet-50 architecture, strives to generate images that can convincingly mislead the discriminator. Through this collaborative interplay, wherein the generator and discriminator are thoughtfully adapted, the GAN framework becomes adept at learning to generate MRI images of exceptional quality. These generated images exhibit a remarkable likeness to authentic MRI scans, consequently elevating the overall standard of image reconstruction quality.

Fig. 3. Discriminator architecture with ResNet50.

*3) Adversarial loss function:* Although there is an empirical rationale for the success of GAN, as mentioned in the introductory section, in practice GAN experiences non-convergence and training fluctuations. The training of the model with loss function, which we shall employ in our model, is said to solve these two issues. It is suggested to minimize the distance between the generator and discriminator as opposed to minimizing JSD. A more stable training history is said to be achieved by this (3):

$$\mathcal{L}(\theta_{G_r}, \theta_{D_r}) = min_{\theta_{G_r}} max_{\theta_{D_r}} \mathbb{E}_{i \sim d}(i)[D_r(i)] - \mathbb{E}_{n \sim d_n(n)}[D_r(G_r(n))]$$
(3)

If the discriminator function $D_r$ is 1-Lipschitz, then this loss function minimizes the distance (4).

$$\|D_r\|_l \leq 1$$
(4)

Weight clipping may be employed to enforce this 1-Lipschitz requirement. The distinct loss functions for the discriminator and generator, which is referred to as $\mathcal{L}_{D_r}$ and $\mathcal{L}_{G_r}$ in (5) and (6) accordingly, for purposes of simplification (3):

$$\mathcal{L}_{D_r}(\theta_{D_r}) = max_{\theta_{D_r}} \mathbb{E}_{i \sim d}(i)[D_r(i)] - \mathbb{E}_{n \sim d_n(n)}[D_r(G_r(n))]$$
(5)

$$\mathcal{L}_A(\theta_{G_r}) = min_{\theta_{G_r}} - \mathbb{E}_{n \sim d_n(n)}[D_r(G_r(n))]$$
(6)

The frequency domain loss ($\mathcal{L}_F$) and normalized root mean square error (NMSE) loss ($\mathcal{L}_N$) to the overall generator loss in addition to $\mathcal{L}_A$ is increased. In the spatial and temporal domains, accordingly, $\mathcal{L}_N$ and $\mathcal{L}_F$ are anticipated to reduce the distinction among fully sampled and reconstructed image. They are characterized as (7) and (8):

$$\mathcal{L}_N(\theta_{G_r}) = \sqrt{\frac{\|i - G_r(n)\|^2}{\|i\|^2}}$$
(7)

$$\mathcal{L}_F(\theta_{G_r}) = \|F\{i\} - F\{G_r(n)\}\|^2 \qquad (8)$$

Here F is the Fourier transform operation to convert an image to a frequency representation, often known as a k-space representation. The generator's final loss function in the simulation is (9):

$$\mathcal{L}_{G_r}(\theta_{G_r}) = \mathcal{L}_A + a\mathcal{L}_N + b\mathcal{L}_F \qquad (9)$$

Here a and b are weights to balance out how much each of the three loss components contributes to the overall loss.

## V. RESULTS AND DISCUSSIONS

Deep learning approaches have been shown to offer significant potential for the development of MRI image reconstruction. Utilizing a wide range of various approaches, study in this field has skyrocketed over the last five years. More thorough evaluation and a broad display of proactively obtained clinical information are required to boost trust in such methods. The work proposes a novel GAN-based MRI reconstruction method that can successfully do away with the input size constraints of GAN. Using ResNet50 connections and loss functions, feature learning can be accelerated. While still working well on images of different sizes, the model has demonstrated improved reconstructing accuracy when compared to previous MRI techniques.

### A. Performance Metrics

In the study focusing on the "Enhancement of MRI Image Reconstruction through Integration of Generative Adversarial Networks (GANs) and Deep Residual Networks (ResNets) in a Hybrid Deep Learning Framework," the evaluation of the proposed methodology is conducted through a range of performance metrics. These metrics include the Peak Signal-to-Noise Ratio (PSNR), measuring image quality by assessing the similarity between original and reconstructed images; the Structural Similarity Index (SSIM) to evaluate the preservation of structural information; Mean Square Error (MSE) and Root Mean Square Error (RMSE) to gauge the average squared differences between original and reconstructed images; Dice Coefficient for segmentation tasks, assessing the overlap in segmented regions; and considerations of computational efficiency, such as inference time, vital for real-time medical imaging applications. These metrics collectively serve as a comprehensive means to quantify the effectiveness and quality of MRI image reconstruction in the context of the hybrid deep learning framework, ensuring it meets the stringent requirements of medical imaging while optimizing computational efficiency.

- Peak Signal-to-Noise Ratio (PSNR): PSNR serves as a metric for assessing the fidelity of the reconstructed MRI images by quantifying the degree of similarity between the original and reconstructed images. Elevated PSNR values are indicative of superior image quality.

- Structural Similarity Index (SSIM): SSIM evaluates the structural information preservation in the reconstructed images. Higher SSIM values indicate that the fine details and structures in the images are retained.

- Mean Square Error (MSE): MSE measures the average squared differences between the original and reconstructed images. Lower MSE values imply better reconstruction quality.

- Root Mean Square Error (RMSE): RMSE is the square root of MSE and provides a more interpretable measure of the reconstruction error.

- Dice Coefficient (for segmentation tasks): If the MRI images are used for segmentation tasks, the Dice coefficient measures the overlap between the segmented regions in the original and reconstructed images.

- Computational Efficiency: This can include metrics like the inference time required for image reconstruction, which is crucial in real-time medical imaging applications.

These performance metrics are essential for quantitatively assessing the quality and efficacy of the MRI image reconstruction using the proposed hybrid deep learning framework, ensuring that it meets the requirements for accurate, high-quality medical imaging while considering computational efficiency.

Outcomes of several reconstruction techniques for a 256×256-pixel image are given. The term "GAN-Resnet50" which refers to the original with the loss function included. The whole set of tests demonstrates that the ResNet50-based generator and discriminator further enhance the model effectiveness whereas the addition of the loss function just minimally accomplishes it.

TABLE I.        Outcome of 256*256 Pixel Images in 10%

| Approaches | Sampling Ratio of 10% | |
|---|---|---|
| | *PSNR* | *SSIM* |
| WDAGAN | 31.689 | 0.920 |
| DAGAN | 30.272 | 0.891 |
| GAN-ResNet50 | 35.281 | 0.950 |

In Table I, the outcomes of different image enhancement approaches for 256x256 pixel images under a 10% sampling ratio are presented, along with their corresponding Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) values. Among the methods evaluated, GAN-ResNet50 demonstrated the highest PSNR at 35.281 and the highest SSIM at 0.950, indicating its superior ability to reconstruct high-quality images from the sparse data, which is crucial for preserving image fidelity. WDAGAN followed with a PSNR of 31.689 and an SSIM of 0.920, showcasing its competitive performance. DAGAN, on the other hand, achieved a PSNR of 30.272 and an SSIM of 0.891, indicating slightly lower image quality compared to the other methods. These results highlight the effectiveness of GAN-ResNet50 in enhancing image reconstruction quality at a 10% sampling ratio, emphasizing its potential for improving image-based diagnostic and analysis tasks in various fields such as medical imaging and remote sensing.

TABLE II.     OUTCOME OF 256×256 PIXEL IMAGES IN 50%

| Approaches | Sampling Ratio of 50% | |
|---|---|---|
| | *PSNR* | *SSIM* |
| WDAGAN | 39.976 | 0.990 |
| DAGAN | 44.939 | 0.954 |
| GAN-ResNet50 | 50.111 | 0.998 |

In Table II, the results for different image enhancement approaches applied to 256×256-pixel images at a 50% sampling ratio are presented, along with their respective Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) values. Notably, GAN-ResNet50 exhibited the highest PSNR at 50.111 and an outstanding SSIM of 0.998, indicating its remarkable ability to reconstruct images with exceptional quality, particularly in scenarios where more data is available. WDAGAN, while also performing well, achieved a PSNR of 39.976 and a very high SSIM of 0.990, highlighting its competence in preserving image details. DAGAN, on the other hand, excelled with an impressive PSNR of 44.939, demonstrating its proficiency in achieving high fidelity image reconstruction even at this increased sampling ratio, although its SSIM of 0.954 suggests some room for improvement in structural similarity. These results underscore the exceptional performance of GAN-ResNet50 in enhancing image quality at a 50% sampling ratio, emphasizing its potential for various image-centric applications where data availability is more substantial, such as high-resolution medical imaging and remote sensing tasks. Fig. 4 shows the Sampling ratio of 256×256-pixel images in 10% and 50%.

The quality of this enhancement is clearly obvious. The internal carotid artery is a bright, vertically directed line that is located slightly below the brain when compared to the equivalent region in the reconstructed image. The reconstruction inaccuracy of this artery in the 'Difference Image' in GAN-ResNet50 is significant. The error looks to be less organized and more random. This suggests that GAN-ResNet50 is performing more qualitatively. The reconstruction error of DAGAN is almost non-existent in the magnified perspective of the red box with corpus callosum and white matter tract and the green box with medial temporal cortex. On the other hand, the image created by DAGAN shows an increased error.

In Table III, the results for various image enhancement approaches applied to 64×128 pixel images at a 10% sampling ratio are displayed, along with their corresponding Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) values. Among the methods assessed, GAN-ResNet50 achieved the highest PSNR at 27.8 and an impressive SSIM of 78.9, demonstrating its capability to enhance image reconstruction quality even when data is sparser. WDAGAN, although with lower PSNR (24.90) and SSIM (76.5) values, exhibited competitive performance. DAGAN, on the other hand, had a PSNR of 24.22 and an SSIM of 74.4, indicating relatively lower image quality compared to the other methods at this low sampling ratio. In Table IV, the results for the same approaches are presented but for 64x128 pixel images with a 50% sampling ratio. GAN-ResNet50 once again displayed the highest PSNR at 40.4 and an outstanding SSIM of 98.1,

signifying its excellent performance in reconstructing high-quality images from more abundant data. WDAGAN and DAGAN also showed substantial improvements at the increased sampling ratio, with competitive PSNR and SSIM values. These results underscore the effectiveness of GAN-ResNet50, especially in scenarios with limited data, and its ability to significantly improve image reconstruction quality, making it a promising solution for various applications, including low-resolution medical imaging and remote sensing tasks. Fig. 5 shows the Sampling ratio of 64×128-pixel images in 10% and 50%.



Fig. 4.   Sampling ratio of 256×256-pixel images in 10% and 50%.



Fig. 5.   Sampling ratio of 64×128 pixel images in 10% and 50%.

TABLE III.    OUTCOME OF 64×128 PIXEL IMAGES IN 10%

| Approaches | Sampling Ratio of 10% | |
|---|---|---|
| | *PSNR* | *SSIM* |
| WDAGAN | 24.90 | 76.5 |
| DAGAN | 24.22 | 74.4 |
| GAN-ResNet50 | 27.8 | 78.9 |

TABLE IV.    OUTCOME OF 64×128 PIXEL IMAGES IN 50%

| Approaches | Sampling Ratio of 50% | |
|---|---|---|
| | *PSNR* | *SSIM* |
| WDAGAN | 38.4 | 97.5 |
| DAGAN | 38.1 | 97.3 |
| GAN-ResNet50 | 40.4 | 98.1 |

This characteristic expands the variety of training samples by enabling the identical GAN-ResNet50 model to analyze an assortment of image dimensions that could be encountered in MR reconstruction. The study has introduced global average pooling to relax the input size requirement of the proposed GAN-ResNet50 framework to provide a fair comparison.

TABLE V.    OUTCOME OF 128×256 PIXEL IMAGES IN 10%

| Approaches | Sampling Ratio of 10% | |
|---|---|---|
| | *PSNR* | *SSIM* |
| WDAGAN | 30.1 | 91.5 |
| DAGAN | 28.3 | 88.6 |
| GAN-ResNet50 | 37.4 | 95.3 |

Table V presents the outcomes for different image enhancement approaches applied to 128x256 pixel images at a 10% sampling ratio, accompanied by their respective Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) values. Notably, GAN-ResNet50 exhibited the highest PSNR at 37.4 and an impressive SSIM of 95.3, signifying its exceptional proficiency in enhancing image reconstruction quality even in situations with limited data. WDAGAN also performed well with a PSNR of 30.1 and a commendable SSIM of 91.5, demonstrating its competence in preserving image details. On the other hand, DAGAN, while achieving a PSNR of 28.3, had a relatively lower SSIM of 88.6, suggesting some room for improvement in structural similarity. These results underscore the superior performance of GAN-ResNet50 in enhancing image quality at a 10% sampling ratio, highlighting its potential for various applications, particularly those involving low-resolution medical imaging and remote sensing, where data availability may be constrained.

TABLE VI.    OUTCOME OF 128×256 PIXEL IMAGES IN 50%

| Approaches | Sampling Ratio of 50% | |
|---|---|---|
| | *PSNR* | *SSIM* |
| WDAGAN | 43.1 | 98.1 |
| DAGAN | 39.7 | 96.9 |
| GAN-ResNet50 | 50.3 | 99.1 |

In Table VI, the results for various image enhancement approaches applied to 128×256-pixel images at a 50% sampling ratio are presented, along with their corresponding Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) values. Remarkably, GAN-ResNet50 demonstrated the highest PSNR at 50.3 and an exceptional SSIM of 99.1, underscoring its remarkable capacity to substantially improve image reconstruction quality in scenarios where data availability is more abundant. WDAGAN also displayed impressive performance, with a PSNR of 43.1 and an outstanding SSIM of 98.1, indicating its competence in preserving image details even at the higher sampling ratio. DAGAN, while achieving a PSNR of 39.7, exhibited a slightly lower SSIM of 96.9, signifying a relatively lower level of structural similarity. These results highlight the outstanding performance of GAN-ResNet50, especially in situations with more data, further emphasizing its potential for diverse applications, including high-resolution medical imaging and remote sensing, where data quality and fidelity are of utmost importance. Fig. 6 shows the Sampling ratio of 128*256 pixel images in 10% and 50%.

In the context of image enhancement and reconstruction, various models have been compared based on their performance metrics, namely, the Structural Similarity Index (SSIM) and Peak Signal-to-Noise Ratio (PSNR).

In Table VII, the study conducted by Hyun et al. [25], the U-net model achieved an SSIM of 0.903, but the PSNR was not specified. Cole et al. [26] implemented an Unsupervised GAN and reported an SSIM of 0.81 along with a PSNR of 26.39. Z. Wang et al. [27] used the Unet-DSSIM model, which achieved an SSIM of 0.88 and a PSNR of 29. In contrast, the proposed method, a Hybrid GAN-DRN model, outperformed the others with an impressive SSIM of 0.99 and a remarkable PSNR of 50.3. These performance metrics highlight the significant advancements achieved by the proposed method in terms of image quality and fidelity. Fig. 7 provides a visual representation of the performance comparison, further emphasizing the superiority of the Hybrid GAN-DRN model in enhancing image reconstruction quality.



Fig. 6.    Sampling ratio of 128×256-pixel images in 10% and 50%.

TABLE VII.    PERFORMANCE COMPARISON

| Reference | Model | SSIM | PSNR |
|---|---|---|---|
| Hyun et al [25] | U-net | 0.903 | - |
| Cole et al. [26] | Unsupervised GAN | 0.81 | 26.39 |
| Z. Wang et al.[27] | Unet-DSSIM | 0.88 | 29 |
| Proposed Method | Hybrid GAN-DRN | 0.99 | 50.3 |



Fig. 7.    Performance comparison of different methods.

The overall performance comparison among the various image enhancement methods clearly demonstrates the remarkable superiority of the proposed Hybrid GAN-DRN model. Its exceptional SSIM score of 0.99 and a remarkable PSNR value of 50.3 signify a substantial leap in image quality and fidelity. The Hybrid GAN-DRN's outstanding SSIM score indicates a near-perfect similarity between the enhanced images and the ground truth, suggesting that it effectively preserves fine image details and structural content. Furthermore, the exceptionally high PSNR value reflects a minimal level of noise and distortion in the reconstructed images, making them highly faithful to the original data. In contrast, the other methods in the comparison, such as U-net, Unsupervised GAN, and Unet-DSSIM, displayed relatively lower SSIM and PSNR values, indicating comparatively reduced image quality and fidelity. This stark performance differential emphasizes the substantial advancements achieved by the proposed Hybrid GAN-DRN model. It excels in enhancing image reconstruction quality, reducing noise, and retaining essential image features. As such, the proposed method holds great promise for a wide array of applications, particularly in fields like medical imaging and remote sensing, where image quality is of paramount importance and where it can significantly improve the accuracy of diagnostic and analysis tasks.

*B.  Discussion*

The study of picture improvement using DRNs and GANs represents a significant advancement in reconstructed images. The study's conclusions and comparison of several methods for enhancing images offer light on the capacity of such models and demonstrate the superior results of the suggested Hybrid GAN-DRN methodology. The primary goal of the study is to compare several image improvement techniques utilizing the SSIM and PSNR, two important performance indicators. Particularly in areas like remote sensing and healthcare imaging, these measures are essential for evaluating the authenticity and quality of augmented images. The analysis of the information shows a distinct pattern: The suggested Hybrid GAN-DRN architecture works noticeably better than each of the approaches compared. This simulation attains a remarkable 50.3 PSNR and a remarkable 0.99 SSIM score. These measurements show how well the suggested approach preserves picture details, lowers noise, and improves the overall level of reconstruction of images. This is especially significant since it creates new opportunities for a variety of image-centric systems in which the precision and integrity of the reconstructions are critical. The practical implications of this research are significant, particularly in areas where diagnostics and choice-making processes are directly impacted by the quality of images. The suggested approach can lead to improved reconstruction of image quality, more precise medical diagnosis, and improved evaluation of information from remote sensing, and improved image-based systems for making decisions. To completely realize the possibilities of the suggested strategy, despite its exceptional contributions, it is imperative to highlight the necessity for additional studies and verification in medical and real-world environments. The study also emphasizes how crucial it is to comprehend and maximize the processing requirements of sophisticated models based on deep learning, particularly for large-scale or instantaneous applications. To sum up, the study of picture improvement using a combination of GANs and DRNs represents a substantial development in the area of reconstruction of images. The exceptional performance of the suggested Hybrid GAN-DRN approach, as shown by SSIM and PSNR parameters, places it in an encouraging spot for enhancing image authenticity as well as quality in a variety of uses, eventually leading to improved accuracy and confidence in image-based decision-making procedures. The benefit of DRGAN is its novel usage of a residual GAN structure to generate a "residual PET map" (RPM), which has the potential to increase PET picture quality with enhanced contextual data. However, more attention on clinical application is needed. FDA-CNN excels in removing MRI artifacts with a completely dense focus, demonstrating promising results. However, the technology requires further investigation to overcome scalability issues to assess its practical application in real-time applications.

VI.  CONCLUSION AND FUTURE WORK

Ultimately, this research presents a potentially effective MRI reconstruction technique based on the GAN-ResNet50 design, tackling the crucial problem of accurate and high-fidelity reconstruction of images. The study's comparison studies show that the suggested strategy outperforms both previous approaches and cutting-edge deep learning methods for maintaining fine-grained image information. This development has significant potential to improve the fidelity and quality of MRI reconstructions, especially in the field of medical imaging because accurate diagnosis and decision-making depend on high-quality images. The acquired findings demonstrate the GAN-ResNet50 strategy's outstanding

performance and its capacity for balancing image texture, border specifics, and total fidelity. Together with the GAN's discriminatory powers, the hybrid deep learning method strikes a compromise that guarantees a degree of reconstructing grade that is both visually realistic and for diagnosis. The suggested GAN-DRN-based MRI reconstruction approach may encounter difficulties in dealing with distinct anatomical variances and diseases, thus limiting generalisation across different clinical circumstances.

*A. Future Work*

There are a number of fascinating avenues for additional study and advancement as researchers look to the years to come. These involve leveraging the GAN-ResNet50 design to optimize real-time MRI reconstruction, verifying the method's clinical effectiveness through rigorous examinations, expanding its uses to multi-modal MRI combination for thorough assessment of patients, and investigating ways to improve robustness, interpretation, ease of use, and transferable learning abilities via explainability methods. These actions are essential for increasing the novel approach's medical acceptance and interpretation, which will ultimately help MRI imaging patients and healthcare providers.

REFERENCES

[1] S. S. Chandra, M. Bran Lorenzana, X. Liu, S. Liu, S. Bollmann, and S. Crozier, "Deep learning in magnetic resonance image reconstruction," J Med Imag Rad Onc, vol. 65, no. 5, pp. 564–577, Aug. 2021, doi: 10.1111/1754-9485.13276.

[2] J. Kim, J. Hong, and H. Park, "Prospects of deep learning for medical imaging," Precis Future Med, vol. 2, no. 2, pp. 37–52, Jun. 2018, doi: 10.23838/pfm.2018.00030.

[3] K. Suzuki, "Overview of deep learning in medical imaging," Radiol Phys Technol, vol. 10, no. 3, pp. 257–273, Sep. 2017, doi: 10.1007/s12194-017-0406-5.

[4] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," Nature, vol. 521, no. 7553, pp. 436–444, May 2015, doi: 10.1038/nature14539.

[5] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting Unreasonable Effectiveness of Data in Deep Learning Era," presented at the Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 843–852. Accessed: Jul. 07, 2023. [Online]. Available: https://openaccess.thecvf.com/content_iccv_2017/html/Sun_Revisiting_Unreasonable_Effectiveness_ICCV_2017_paper.html

[6] A. Elhadad, F. Alanazi, A. I. Taloba, and A. Abozeid, "Fog Computing Service in the Healthcare Monitoring System for Managing the Real-Time Notification," Journal of Healthcare Engineering, vol. 2022, pp. 1–11, Mar. 2022, doi: 10.1155/2022/5337733.

[7] J. Schmidhuber, "Deep learning in neural networks: An overview," Neural Networks, vol. 61, pp. 85–117, Jan. 2015, doi: 10.1016/j.neunet.2014.09.003.

[8] A. Taloba, M. A. Fouly, and T. Soliman, "Developing an Efficient Secure Query Processing Algorithm on Encrypted Databases using Data Compression," 2022.

[9] S. Pouyanfar et al., "A Survey on Deep Learning: Algorithms, Techniques, and Applications," ACM Comput. Surv., vol. 51, no. 5, pp. 1–36, Sep. 2019, doi: 10.1145/3234150.

[10] F. Xing, Y. Xie, H. Su, F. Liu, and L. Yang, "Deep Learning in Microscopy Image Analysis: A Survey," IEEE Trans. Neural Netw. Learning Syst., vol. 29, no. 10, pp. 4550–4568, Oct. 2018, doi: 10.1109/TNNLS.2017.2766168.

[11] M. Bakator and D. Radosav, "Deep Learning and Medical Diagnosis: A Review of Literature," MTI, vol. 2, no. 3, p. 47, Aug. 2018, doi: 10.3390/mti2030047.

[12] M. Caramia and E. Pizzari, "A Bi-objective cap-and-trade model for minimising environmental impact in closed-loop supply chains," Supply Chain Analytics, vol. 3, p. 100020, Sep. 2023, doi: 10.1016/j.sca.2023.100020.

[13] M. O. Malavé et al., "Reconstruction of undersampled 3D non-Cartesian image-based navigators for coronary MRA using an unrolled deep learning model," Magn Reson Med, vol. 84, no. 2, pp. 800–812, Aug. 2020, doi: 10.1002/mrm.28177.

[14] S. Ramani, Zhihao Liu, J. Rosen, J. Nielsen, and J. A. Fessler, "Regularization Parameter Selection for Nonlinear Iterative Image Restoration and MRI Reconstruction Using GCV and SURE-Based Methods," IEEE Trans. on Image Process., vol. 21, no. 8, pp. 3659–3672, Aug. 2012, doi: 10.1109/TIP.2012.2195015.

[15] Q. Du, Y. Qiang, W. Yang, Y. Wang, Y. Ma, and M. B. Zia, "DRGAN: a deep residual generative adversarial network for PET image reconstruction," IET Image Processing, vol. 14, no. 9, pp. 1690–1700, Jul. 2020, doi: 10.1049/iet-ipr.2019.1107.

[16] Md. B. Hossain, K.-C. Kwon, S. M. Imtiaz, O.-S. Nam, S.-H. Jeon, and N. Kim, "De-Aliasing and Accelerated Sparse Magnetic Resonance Image Reconstruction Using Fully Dense CNN with Attention Gates," Bioengineering, vol. 10, no. 1, p. 22, Dec. 2022, doi: 10.3390/bioengineering10010022.

[17] M. V. R. Manimala, C. Dhanunjaya Naidu, and M. N. Giri Prasad, "Sparse MR Image Reconstruction Considering Rician Noise Models: A CNN Approach," Wireless Pers Commun, vol. 116, no. 1, pp. 491–511, Jan. 2021, doi: 10.1007/s11277-020-07725-0.

[18] R. Hou and F. Li, "IDPCNN: Iterative denoising and projecting CNN for MRI reconstruction," Journal of Computational and Applied Mathematics, vol. 406, p. 113973, May 2022, doi: 10.1016/j.cam.2021.113973.

[19] S. V. Yakkundi and D. P. Subha, "Convolutional LSTM: A Deep learning approach for Dynamic MRI Reconstruction," in 2020 4th International Conference on Trends in Electronics and Informatics (ICOEI)(48184), Tirunelveli, India: IEEE, Jun. 2020, pp. 1011–1015. doi: 10.1109/ICOEI48184.2020.9142982.

[20] S. Huang, L. Sun, M. Yousefnezhad, M. Wang, and D. Zhang, "Temporal Information-Guided Generative Adversarial Networks for Stimuli Image Reconstruction From Human Brain Activities," IEEE Trans. Cogn. Dev. Syst., vol. 14, no. 3, pp. 1104–1118, Sep. 2022, doi: 10.1109/TCDS.2021.3098743.

[21] Y. Wu, Y. Ma, J. Liu, J. Du, and L. Xing, "Self-attention convolutional neural network for improved MR image reconstruction," Information Sciences, vol. 490, pp. 317–328, Jul. 2019, doi: 10.1016/j.ins.2019.03.080.

[22] S. Wang et al., "DeepcomplexMRI: Exploiting deep residual network for fast parallel MR imaging with complex convolution," Magnetic Resonance Imaging, vol. 68, pp. 136–147, May 2020, doi: 10.1016/j.mri.2020.02.002.

[23] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath, "Generative adversarial networks: An overview," IEEE signal processing magazine, vol. 35, no. 1, pp. 53–65, 2018.

[24] C. Xu, J. Tao, Z. Ye, J. Xu, and W. Kainat, "Adversarial training and dilated convolutions for compressed sensing MRI," in Eleventh International Conference on Digital Image Processing (ICDIP 2019), SPIE, 2019, pp. 1014–1021.

[25] C. M. Hyun, H. P. Kim, S. M. Lee, S. Lee, and J. K. Seo, "Deep learning for undersampled MRI reconstruction," Physics in Medicine & Biology, vol. 63, no. 13, p. 135007, 2018.

[26] E. K. Cole, J. M. Pauly, S. S. Vasanawala, and F. Ong, "Unsupervised MRI reconstruction with generative adversarial networks," arXiv preprint arXiv:2008.13065, 2020.

[27] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," IEEE transactions on image processing, vol. 13, no. 4, pp. 600–612, 2004.

# Hybrid Vision Transformers and CNNs for Enhanced Transmission Line Segmentation in Aerial Images

Hoanh Nguyen*, Tuan Anh Nguyen

Faculty of Electrical Engineering Technology, Industrial University of Ho Chi Minh City, Ho Chi Minh City, Vietnam

*Abstract*—This paper presents a novel architecture for the segmentation of transmission lines in aerial images, utilizing a hybrid model that combines the strengths of Vision Transformers (ViTs) and Convolutional Neural Networks (CNNs). The proposed method first employs a Swin Transformer backbone (Swin-B) that processes the input image through a hierarchical structure, effectively capturing multi-scale contextual information. Following this, an upsampling strategy is employed, wherein the features extracted by the transformer are refined through convolutional layers, ensuring that the resolution is maintained, and spatial details are recovered. To integrate multi-level feature maps, a feature fusion module with a squeeze-and-excitation (SE) layer is introduced, which consolidates the benefits of both high-level and low-level feature extractions. The SE layer plays a pivotal role in augmenting the feature channels, focusing the model's attention on the most informative features for transmission line detection. By leveraging the global receptive field of ViTs for comprehensive context and the local precision of CNNs for fine-grained detail, our method aims to set a new benchmark for transmission line segmentation in aerial imagery. The effectiveness of our approach is demonstrated through extensive experiments and comparisons with existing state-of-the-art methods.

*Keywords—Vision transformers; convolutional neural networks; transmission lines segmentation; hybrid model; feature fusion*

## I. INTRODUCTION

Transmission line segmentation in aerial images is a critical task in the maintenance and monitoring of electrical power grids. It enables the automated inspection of power lines for fault detection, vegetation encroachment, and structural analysis, which are essential for ensuring the reliability and safety of electricity distribution. However, this task is fraught with challenges. Aerial images often have highly variable lighting conditions, weather effects, and diverse landscapes that can obscure the visibility of transmission lines. Additionally, the lines themselves can be difficult to distinguish due to their thin and linear nature against complex backgrounds. Another significant challenge is the presence of other linear structures, such as roads and railways, that can be easily confused with power lines by automated systems. The movement of the aerial platform, whether it's a drone or a manned aircraft, introduces motion blur and varying angles of capture, further complicating the segmentation process. Addressing these challenges requires robust algorithms capable of high precision and adaptability to a range of environmental conditions and image qualities. Traditional methods for vision-based transmission line detection and segmentation have

revolved around the utilization of edge and line segment detection techniques as foundational steps. These methods generate a multitude of potential cable segment candidates by applying algorithms such as the Hough transform, Radon transform, and various other heuristic and search-based line detection strategies, such as the circle-based search, heuristic line detection, Line Segment Detector (LSD), and Edge Drawing for line segment detection (EDLines). Once edges and line segments are detected, a set of specialized rules, informed by the structural characteristics of cables and the context of their surroundings, are applied to discern correct cable segments and eliminate false positives. One of the main drawbacks of these approaches is their dependency on numerous parameters and complex rules that need to be meticulously set by hand, which hinders their adaptability to different environments. As a result, the precision and robustness of traditional methods can be significantly compromised due to environmental variations, making it challenging to maintain consistent performance across diverse scenarios. The advent of deep learning has catalyzed significant advancements in the domain of vision-based transmission line detection and segmentation. CNN-based methods have eclipsed traditional techniques, demonstrating substantial improvements in detection accuracy and computational efficiency. These deep neural networks facilitate end-to-end learning and inference, simplifying the complex parameter tuning process inherent in multistage approaches and enhancing generalizability across varied scenarios. For instance, CNNs have been trained to identify image patches containing cables, which are then further processed using traditional methods like the Hough transform for line segmentation. Moreover, some approaches have integrated fully convolutional networks with line segment regressors for direct line segment detection, particularly effective in scenarios where aerial images capture transmission lines at close range. However, when dealing with long-range and wide-angle captures that result in indistinct or slightly curved cable representations, these methods pivot towards a pixel-wise segmentation framework, employing semantic and instance segmentation techniques to provide a more nuanced cable detection and enable individual cable instance identification, which is pivotal for autonomous UAV applications.

Recent years have witnessed the rapid development of ViTs [1, 2]. ViTs leverages the transformer architecture, originally designed for natural language processing, to handle sequences of image patches as input. ViTs model relationships between these patches through self-attention mechanisms, making them capable of capturing global dependencies within an image. The

combination of CNNs and ViT into a hybrid architecture aims to harness the local feature extraction proficiency of CNNs with the global context understanding of ViTs. CNNs are adept at recognizing patterns and textures within small regions of an image, making them excellent for tasks that require detailed local information, such as edge detection. ViTs, with their attention-based approach, can consider the entire image at once, which allows for a more holistic understanding of the scene. By integrating both, the hybrid model can effectively process and integrate both local and global information, leading to improved performance on complex tasks like transmission line segmentation. This synergy can provide a more nuanced understanding of images, enabling the model to be both precise in detail and comprehensive in scope, potentially overcoming limitations found in models that rely on a single approach.

This study introduces an innovative hybrid architecture for the precise segmentation of transmission lines in aerial images, leveraging the synergistic potentials of ViTs and CNNs. The core of our proposed method is a Swin Transformer backbone, adept at hierarchically processing the input image to encapsulate multi-scale contextual information. This is complemented by an upsampling mechanism that meticulously refines the transformer-extracted features via convolutional layers, crucial for preserving resolution and restoring spatial details. A feature fusion module, equipped with an SE layer, is integrated to merge feature maps from multiple levels, harnessing both the high-level and low-level extraction strengths. The SE layer is instrumental in enhancing feature channels, directing the model's focus towards the most salient features for detecting transmission lines. Our approach is designed to exploit the expansive receptive field of ViTs for global context awareness, while utilizing the CNNs' local precision for capturing intricate details, thereby establishing a new standard for transmission line segmentation in aerial photography. The method's superiority is validated through comprehensive experimental benchmarks, showcasing its advancement over current leading methodologies.

The rest of the paper is organized as follows: Section II presents related studies; Section III details our proposed model; Section IV describes the experiments and results; Section V provides the conclusions.

## II. RELATED WORK

### A. Transmission Line Detection and Segmentation

Traditional approaches to vision-based detection and segmentation of transmission lines have primarily focused on employing edge and line segment detection techniques as their fundamental processes. In study [3], a real-time algorithm was developed for detecting power lines in UAV video images, where the process begins with converting video images into binary images using adaptive thresholding. Subsequently, Hough Transform identifies line candidates in these binary images, and a fuzzy C-means clustering algorithm discriminates actual power lines from these candidates. Mu et al. [4] proposed a method for automatically extracting power lines from cluttered natural backgrounds in aerial images. The approach involves using a Gabor filter to eliminate background noise, followed by the application of the Hough transform to detect straight lines in the images. Zhang et al. [5] introduced a new method for detecting and tracking power lines, starting with the use of the Hough transform to extract line segments. The method then employs K-means clustering in the Hough space to filter and identify power lines and utilizes a Kalman filter for tracking these lines within the continuity of a video sequence. In study [6], the authors presented an algorithm that capitalizes on the geometric relationships inherent to circle symmetry for line segment detection. It employs Canny and Steerable Filters to detect line segments, which are then linked in a subsequent stage for effective analysis. Sharma et al. [7] introduced a novel morphological operator and robust image space heuristics for the accurate location and complete extraction of power lines. Santos et al. [8] introduced PLineD, a new vision-based power line detection algorithm designed to robustly detect power lines, even in noisy image backgrounds. Although traditional approaches to vision-based detection and segmentation of transmission lines have achieved some success, they still have many limitations. A significant limitation of these methods is their reliance on a multitude of parameters and intricate rules that require careful manual adjustment, impeding their flexibility across various environments. Consequently, the accuracy and reliability of traditional approaches are often adversely affected by environmental changes, posing challenges in achieving uniform performance in different settings. With the outstanding advantages of CNNs, many methods using CNNs for transmission line segmentation have been proposed [9]. In [10], the authors introduced a pyramidal patch classification framework that effectively eliminates clutter without relying on additional auxiliary tools. This is achieved through a hierarchical patch partition and selection strategy, complemented by a new spatial grid pooling layer in the CNN-based classifier. Nguyen et al. [11] presented LS-Net, a rapid, single-shot line-segment detector tailored for power line detection, which is fully convolutional by design and comprises three modules: a fully convolutional feature extractor, a classifier, and a line segment regressor. Lee et al. [12] presented a weakly supervised learning algorithm for identifying power lines. The algorithm classifies sub-regions within images using a sliding window approach and a CNN. In [13], a Transmission Line Detection (TLD) algorithm, CableNet, is proposed, drawing inspiration from instance segmentation and incorporating enhancements to Fully Convolutional Networks (FCNs) [14] with overlaying dilated and spatial convolutional layers for better representation of transmission lines, and dual output branches for generating multidimensional feature maps for instance segmentation.

### B. Vision Transformer-CNN Hybrid Models

In recent years, the rapid development of ViTs has significantly advanced the field of computer vision, leading to their widespread application in tasks ranging from image classification to complex scene understanding [15, 16]. Instead of simplifying ViTs, another prominent research direction involves merging components of ViTs and CNNs to create novel backbone architectures. These hybrid models combine the local feature extraction prowess of CNNs with the global contextual understanding afforded by ViTs, thus offering a comprehensive approach to image analysis. Follow this approach, [17] highlighted the adaptation of principles from the

extensive literature on CNNs, particularly the use of activation maps with decreasing resolutions, to enhance the design of transformers. The study in [18] investigated the optimization challenges of ViT models, attributing the issues to their 'patchify stem' design, and proposes a solution by replacing it with stacked stride-two 3x3 convolutions. This modification significantly enhances optimization stability and model accuracy, leading to the recommendation of using a standard, lightweight convolutional stem in ViT models for improved performance and robustness. In study [19], the authors introduced BoTNet, a versatile and efficient backbone architecture for various computer vision tasks, which enhances performance by integrating self-attention mechanisms. This is achieved by replacing spatial convolutions with global self-attention in the last three bottleneck blocks of a ResNet, leading to notable improvements in instance segmentation and object detection, while simultaneously reducing the number of parameters and maintaining minimal latency overhead. ConViT [20] presented the concept of gated positional self-attention (GPSA), a novel form of positional self-attention designed with a flexible, 'soft' convolutional inductive bias. GPSA layers are initially configured to emulate the locality characteristic of convolutional layers but are also equipped with a gating parameter that allows each attention head to dynamically balance the focus between positional and content information. Guo et al. [21] proposed a novel hybrid network that synergizes the long-range dependency capturing capabilities of transformers with the local information extraction prowess of CNNs. Recently, PVTv1 [22], PVTv2 [23], LITv1 [24], and LITv2 [25] incorporate convolutional operations at each stage of ViT models to diminish the token count and construct hybrid, multi-stage structures.

## III. METHOD

### A. Model Architecture

Fig. 1 illustrates the overall pipeline of our method, which integrates a vision transformer encoder with a convolutional neural network decoder to create a hybrid model for the segmentation of transmission lines in aerial images. Input images undergo a hierarchical processing through multiple layers of the Swin Transformer [26], each reducing the spatial dimensions while increasing the depth of feature representation. These layers (Layer 1 to Layer 4) progressively transform the input, capturing intricate details and contextual information at various scales. The transformed features are then upsampled and passed through convolutional layers to refine the feature maps, ensuring that spatial information is preserved and enhanced. The upsampling process gradually restores the resolution of the feature maps, which are then combined through feature fusion steps. These fusion steps are essential as it aggregates multi-scale information, enabling the model to capture both high-level semantic information and low-resolution spatial details. A squeeze-and-excitation (SE) layer is subsequently employed to recalibrate the feature channels, emphasizing informative features while suppressing less useful ones. Finally, a segmentation head, comprising a series of convolutional layers, is responsible for generating the output segmentation map that delineates the transmission lines within the aerial images.



Fig. 1. Model architecture.

Fig. 2. The structure of Swin Transformer encoder.

TABLE I. DETAILED ARCHITECTURE OF SWIN-B

| Layer | Output Size | Number of Blocks | Attention Heads | Attention Head Dimensions | MLP Ratio |
|---|---|---|---|---|---|
| Patch Partition | w/4 × h/4 × 48 | N/A | N/A | N/A | N/A |
| Layer 1 | w/4 × h/4 × 96 | 2 | 3 | 32 | 4 |
| Layer 2 | w/8 × h/8 × 192 | 2 | 6 | 32 | 4 |
| Layer 3 | w/16 × h/16 × 384 | 6 | 12 | 32 | 4 |
| Layer 4 | w/32 × h/32 × 768 | 2 | 24 | 32 | 4 |

## B. Swin-B-based Encoder

Fig. 2 illustrates the structure of Swin Transformer encoder. The original input image is represented as $w \times h \times 3$, where $w$ and $h$ are the width and height of the image, and $3$ represents the RGB color channels. This image is then partitioned into patches. The size of these patches is 4×4 pixels, which are then flattened and linearly embedded into a higher-dimensional space (e.g., 48 features per patch). So, the input dimension to the first transformer layer is $w/4 \times h/4 \times 48$. As the input passes through each Swin Transformer layer consisting consecutive $N_i$ ($i = 1, 2, 3, 4$) Swin transformer blocks, the spatial resolution is further reduced, and the feature dimensionality is increased, thus enhancing the model's ability to capture more complex features at different scales. The Swin Transformer employs a self-attention mechanism within each block. The attention is computed using queries ($Q$), keys ($K$), and values ($V$), which are derived from the input feature maps. The formulation of self-attention within the Swin Transformer involves a sequence of operations beginning with the computation of $Q$, $K$, and $V$ via linear transformations of the input feature map. The attention scores are then determined by calculating the dot product between $Q$ and $K$. These scores are normalized using a softmax function to derive attention weights, which are subsequently used to obtain a weighted feature representation by multiplying them with $V$. To ensure stability in the gradients during training, a scaling factor, commonly the inverse square root of the keys' dimensionality, is optionally applied to the dot product of $Q$ and $K$. Mathematically, the attention can be represented as:

$$Attention(Q, K, V) = Softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

where, $d_k$ is the dimensionality of the key vectors. This allows the model to focus on different parts of the image depending on the learned importance of each feature.

The Swin Transformer framework is offered in four distinct variants, namely Swin-T (Tiny), Swin-S (Small), Swin-B (Base), and Swin-L (Large), each differing in capacity and computational requirements. The choice of a particular Swin Transformer variant for a given task hinges on a balance between the model's empirical performance and the computational constraints of the available hardware. In the case of the segmentation of transmission lines in aerial images, Swin-B was selected due to its robust performance in capturing intricate details and providing a higher feature resolution necessary for the precise delineation of transmission lines, which are often slender and require fine-grained detection capabilities. Moreover, Swin-B strikes a balance between computational efficiency and model complexity, making it a pragmatic choice for tasks demanding high accuracy without exceedingly intensive computational demands. Table I provides detailed architecture of Swin-B backbone.

## C. CNNs-based Decoder

The decoder leverages a U-Net-like structure known for its effectiveness in segmentation tasks due to its ability to combine low-level feature maps with high-level ones, thus capturing context and fine details. Each feature map output from the Swin-B backbone passes through a 3×3 convolution layer. This operation serves to refine the feature maps by applying filters that can capture spatial hierarchies within the data. After the convolution layers, the feature maps are upsampled. This process increases the spatial resolution of the feature maps to prepare them for feature fusion. The upsampling doubles the height and width of the feature maps, as is common in U-Net architectures [27] to match the dimensions of the feature maps from the encoder that will be fused. The upsampled feature maps are then fused with corresponding feature maps from earlier layers of the encoder. This step is crucial as it reintroduces higher resolution details that may have been lost during downsampling in the encoder. Feature fusion is done using element-wise addition operation. The last part of the decoder is the segmentation head, which

outputs the final segmentation map. This head takes the processed feature maps and applies a combination of convolutional layers, activation functions, and sigmoid layer to generate the pixel-wise classification of the transmission lines. Given the size of the input to the decoder is $\frac{H}{32} \times \frac{W}{32} \times 8C$, the size of the output should match the original height and width of the input image $H \times W \times 2$.

### D. Squeeze-and-Excitation Layer

Before the final output, we employ a Squeeze-and-Excitation (SE) layer [28] to recalibrate the feature channels by explicitly modelling the interdependencies between them. The SE layer uses global average pooling to squeeze global spatial information into a channel descriptor, then uses two fully connected layers to capture channel-wise dependencies, and finally applies the channel weights back to the original feature maps to emphasize useful features and suppress less useful ones. Fig. 3 shows the architecture of the SE layer. Let the input to the SE layer be a feature map $F$ with dimensions $H' \times W' \times C$, where $H'$ and $W'$ are the spatial dimensions after upsampling and $C$ is the number of channels. The SE layer first performs a global average pooling operation on $F$, which squeezes the spatial dimensions $H'$ and $W'$ into a single channel descriptor $z$ with dimensions $1 \times 1 \times C$. Mathematically, this is represented as:

$$z_c = \frac{1}{H' \times W'} \sum_{i=1}^{H'} \sum_{j=1}^{W'} F_{i,j,c} \quad (2)$$



Fig. 3. The architecture of the SE layer.

where, $z_c$ is the *c-th* element of $z$; $F_{i,j,c}$ is the value at position *(i, j)* in channel $c$ of the feature map $F$.

The SE layer then passes $z$ through two fully connected layers. The first layer reduces the channel dimensionality from $C$ to $\frac{C}{r}$ using a ReLU activation function, and the second layer increases it back to $C$ using a sigmoid activation function, thus generating the channel-wise weights $s$ with the same dimension $1 \times 1 \times C$. Mathematically, this is represented as:

$$s = \sigma\big(g(z,W)\big) = \sigma(W_2.\theta(W_1.z)) \quad (3)$$

where, $\sigma$ denotes the sigmoid activation, $\theta$ denotes the ReLU activation, $W_1$ and $W_2$ are the weights of the fully connected layers, and $g$ represents the excitation function.

Finally, the SE layer applies these weights $s$ back to the original feature map $F$ through channel-wise multiplication, producing the output feature map $F'$ with the same spatial dimensions $H' \times W'$ but with recalibrated channels:

$$F'_{i,j,c} = s_c.F_{i,j,c} \quad (4)$$

This operation scales each channel of the input feature map by the corresponding learned weight, emphasizing informative features and suppressing less relevant ones. The output of the SE layer is then ready to be passed to the subsequent layers in the decoder for further processing towards the final segmentation map.

### E. Loss Function

Binary Cross-Entropy (BCE) loss is a commonly used loss function for binary classification tasks, such as the segmentation of transmission lines in aerial images, where each pixel is classified as either belonging to a transmission line (positive class) or background (negative class). The BCE loss function measures the distance between the predicted probabilities and the actual binary labels, penalizing predictions that diverge from the true labels. Formally, the BCE loss for a single pixel is calculated as:

$$L_i = -[y log(p) + (1-y)log(1-p)] \quad (5)$$

where, $y$ is the true label of the pixel, and $p$ is the predicted probability that the pixel belongs to the transmission line class. The true label $y$ is 1 if the pixel is part of a transmission line and 0 otherwise. The predicted probability $p$ is obtained from the output of a sigmoid activation function in the last layer of the neural network, ensuring that $p$ is in the range [0,1].

For the entire image, the total BCE loss is the average of the individual pixel losses:

$$L_{BCE} = -\frac{1}{N}\sum_{i=1}^{N} L_i \quad (6)$$

where, $N$ is the total number of pixels in the image.

## IV. EXPERIMENTS

### A. Dataset and Metrics

We use the TTPLA dataset [29] to evaluate the proposed method. The TTPLA dataset is a specialized collection of aerial images designed for the detection and segmentation of transmission towers and power lines. This dataset is significant for training and evaluating machine learning models, particularly in the domain of remote sensing and automated monitoring of electrical infrastructure. It includes high-resolution images that capture the intricate details of transmission towers and power lines from various angles and under different lighting conditions. The diversity of the dataset aids in developing robust models capable of accurately identifying and segmenting these structures. The dataset consists of 1,100 images with a resolution of 3,840×2,160 pixels and manually labeled 8,987 instances of transmission lines and transmission towers. For the purpose of evaluating the transmission line segmentation task, we only employ the labels of transmission lines for training and testing.

Precision (*P*), recall (*R*), Intersection over Union (*IoU*), and *F-score* are critical metrics for evaluating the performance of models in the segmentation of transmission lines in aerial images. Precision measures the ratio of correctly predicted positive observations to the total predicted positives. It is formulated as:

$$P = \frac{TP}{TP+FP} \quad (7)$$

where, *TP* is true positives and *FP* is false positives.

Recall assesses the ratio of correctly predicted positive observations to all actual positives. It's given by:

$$R = \frac{TP}{TP+FN} \qquad (8)$$

with *FN* being false negatives.

Intersection over Union (*IoU*), also known as the Jaccard index, is the area of overlap between the predicted segmentation and the ground truth divided by the area of union. The formula is:

$$IoU = \frac{\text{area of overlap}}{\text{area of union}} \qquad (9)$$

*F-score* is the harmonic mean of precision and recall, providing a balance between them. It's calculated by:

$$F - score = 2 \times \frac{P \times R}{P+R} \qquad (10)$$

These metrics are pivotal for tuning models to the specific challenges of aerial image segmentation, such as delineating thin and often indistinct transmission lines against complex backgrounds. High precision indicates a model that reliably identifies line pixels, while high recall shows it finds most of the actual line pixels. *IoU* gives an overall sense of the model's accuracy, and *F-score* offers a single measure to assess both precision and recall.

### B. Implementation Details

We leverage the strengths of transformer models, specifically building upon the Swin Transformer for the encoder component. The Swin-B model is pretrained on ImageNet-22k with a resolution of 384×384, maintaining the window size (M) as in the pretrained models. To adapt to the higher resolution requirements of the semantic segmentation task, we fine-tune these models based on the dataset's resolution. Following methodologies in the literature, relative position bias is incorporated when calculating attention scores. The decoders are initialized with random weights from a normal distribution. For optimization, the AdamW optimizer [30] is employed. The input aerial images are resized to a standard resolution of 512×512 pixels, striking a balance between preserving detail and maintaining computational efficiency. The learning rate is set at 1e-4, paired with a cosine

annealing scheduler to facilitate adaptive learning rate adjustments over approximately 100 training epochs. The model's training is executed on a high-end NVIDIA RTX 4080 GPU. A batch size of 4 is used. Implementation is carried out using deep learning frameworks PyTorch. To further bolster the model's ability to generalize, data augmentation techniques including random rotations, flipping, scaling, and brightness adjustments are employed.

### C. Comparison with Existing Methods on TTPLA Dataset

We compared the proposed model with six existing models on the TTPLA Dataset, as shown in Table II. It is evident that our model, which integrates a vision transformer encoder with a convolutional neural network decoder, outperforms most of the existing models in terms of precision, IoU (Intersection over Union), and F-score. These metrics are critical for assessing the effectiveness of segmentation models in aerial imagery. Notably, the proposed model achieves a precision of 0.855 and an F-score of 0.671, surpassing the UNet, UNet++, and Focal-UNet models that also use the Resnet-18 architecture. This superior performance can be attributed to the efficient feature representation and fusion enabled by the hybrid architecture of the proposed model. The Swin Transformer layers in the encoder capture intricate details and contextual information at various scales, which is crucial for accurately delineating transmission lines in complex aerial images. The use of the squeeze-and-excitation layer further enhances the model by emphasizing informative features, allowing for a more nuanced segmentation output. This is evident in the comparative improvement in precision and F-score, where the model excels in correctly identifying relevant pixels while maintaining high overall segmentation accuracy. In contrast, models like LCNN and HAWP, based on the Hourglass architecture, show significantly lower precision and F-score values, indicating a lesser ability to accurately segment transmission lines in aerial images. Their lower performance might be due to less effective feature extraction and fusion compared to the proposed hybrid model. Overall, the proposed model's superior performance across multiple metrics, especially in precision and F-score, highlights its effectiveness in segmenting transmission lines in aerial images, demonstrating the advantages of its novel architecture combining vision transformer and convolutional layers.

TABLE II.        SEGMENTATION PERFORMANCE OF THE PROPOSED METHOD AND THE COMPARISON METHODS ON THE TTPLA DATASET

| Models | Backbone | P | R | IoU | F-score |
|---|---|---|---|---|---|
| DeepLabv3+ [31] | Resnet-18 | 0.784 | 0.510 | 0.424 | 0.573 |
| UNet [27] | Resnet-18 | 0.846 | 0.583 | 0.515 | 0.662 |
| UNet++ [32] | Resnet-18 | 0.843 | 0.591 | 0.522 | 0.668 |
| Focal-UNet [33] | Resnet-18 | 0.784 | 0.577 | 0.504 | 0.662 |
| LCNN [34] | Hourglass | 0.541 | 0.315 | 0.498 | 0.519 |
| HAWP [35] | Hourglass | 0.581 | 0.421 | 0.485 | 0.532 |
| Our Model | Swin-B | 0.855 | 0.579 | 0.522 | 0.671 |

Fig. 4.   Sample transmission line segmentation results of the proposed model.

Fig. 4 displays a side-by-side comparison of original aerial images against the segmentation results produced by the proposed model for transmission line identification. Across various landscapes such as residential areas, road intersections, and open fields with transmission towers, the model delineates the transmission lines with a high degree of precision, as indicated by the blue lines overlaying the images. Specifically, the model accurately identifies transmission lines in residential areas without confusion from similar-colored backgrounds, showing its precision in challenging environments. At a road intersection, the model successfully differentiates transmission lines from road markings despite visual noise, indicating strong feature extraction capabilities. In open fields, the model effectively handles contrasts and textures, maintaining accurate segmentation over uniform backgrounds like grass. These results underline the model's robustness in varied settings and its applicability in real-world tasks such as infrastructure monitoring from aerial imagery.

### D. Comparison of Swin-B with Other State-of-the-Art Backbones

We conducted experiments to evaluate the performance of Swin-B encoder. Fig. 5 shows the F-score performance of various backbone architectures employed in image segmentation models. Swin-B tops the chart with an F-score of 0.671, indicating its superior ability in combining features effectively for precise segmentation tasks. Swin-T closely trails with a marginally lower F-score of 0.668, suggesting that while it is slightly less effective than Swin-B, it remains a highly competitive architecture. ResNeSt-101 and ResNet-101, both advanced iterations of the ResNet family, score 0.642 and 0.630 respectively, pointing to a proficient but noticeably lesser segmentation capability compared to the Swin architectures. VGG-16, the oldest architecture among those compared, shows its limitations with an F-score of 0.585, underscoring the advancements in backbone architectures for segmentation tasks and the importance of choosing the right one for optimal performance.

Fig. 5. F-score comparison of different backbone architectures.

## V. CONCLUSION

This paper presents a novel hypbid architecture for the segmentation of transmission lines in aerial images. The presented hybrid segmentation model, which leverages the synergy of a vision transformer encoder and a convolutional neural network decoder, has proven highly effective in the segmentation of transmission lines from aerial images. The model's performance, as demonstrated on the TTPLA Dataset, is superior to existing models, achieving remarkable precision and F-score metrics. Its ability to handle complex backgrounds and maintain high accuracy in diverse environments showcases its robustness and adaptability. The successful application of this model paves the way for its integration into aerial survey systems, offering significant improvements in the monitoring and maintenance of power line infrastructures, potentially reducing costs and increasing operational efficiency. The research outcomes not only contribute to the advancement of segmentation techniques but also underline the transformative impact of integrating transformer architectures within computer vision tasks. For future work, we will incorporate advanced object detection algorithms to identify not just transmission lines but also associated structures such as towers and insulators. This would provide a more comprehensive analysis of the aerial imagery, facilitating detailed inspections and maintenance planning.

## REFERENCES

[1] A. Dosovitskiy *et al.*, "An image is worth 16×16 words: Transformers for image recognition at scale," in *Proc. 9th Int. Conf. Learn. Represent.*, Virtual Event, Austria, May 2021, pp. 1–22.

[2] Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

[3] Yang, Tang Wen, Hang Yin, Qiu Qi Ruan, Jian Da Han, Jun Tong Qi, Qing Yong, Zi Tong Wang, and Zeng Qi Sun. "Overhead power line detection from UAV video images." In *2012 19th International Conference on Mechatronics and Machine Vision in Practice (M2VIP)*, pp. 74-79. IEEE, 2012.

[4] Mu, Chao, Jie Yu, Yanming Feng, and Jinhai Cai. "Power lines extraction from aerial images based on Gabor filter." In *International Symposium on Spatial Analysis, Spatial-Temporal Data Modeling, and Data Mining*, vol. 7492, pp. 1081-1088. SPIE, 2009.

[5] Zhang, Jingjing, Liang Liu, Binhai Wang, Xiguang Chen, Qian Wang, and Tianru Zheng. "High speed automatic power line detection and

tracking for a UAV-based inspection." In *2012 International Conference on Industrial Control and Electronics Engineering*, pp. 266-269. IEEE, 2012.

[6] Cerón, Alexander, and Flavio Prieto. "Power line detection using a circle based search with UAV images." In *2014 international conference on unmanned aircraft systems (ICUAS)*, pp. 632-639. IEEE, 2014.

[7] Sharma, Hrishikesh, Rajeev Bhujade, V. Adithya, and P. Balamuralidhar. "Vision-based detection of power distribution lines in complex remote surroundings." In *2014 Twentieth National Conference on Communications (NCC)*, pp. 1-6. IEEE, 2014.

[8] Santos, Tiago, Miguel Moreira, J. Almeida, André Dias, Alfredo Martins, J. Dinis, J. Formiga, and E. Silva. "PLineD: Vision-based power lines detection for Unmanned Aerial Vehicles." In *2017 IEEE International Conference on Autonomous Robot Systems and Competitions (ICARSC)*, pp. 253-259. IEEE, 2017.

[9] Jenssen, Robert, and Davide Roverso. "Automatic autonomous vision-based power line inspection: A review of current status and the potential role of deep learning." *International Journal of Electrical Power & Energy Systems* 99 (2018): 107-120.

[10] Li, Yan, Chaofeng Pan, Xianbin Cao, and Dapeng Wu. "Power line detection by pyramidal patch classification." *IEEE Transactions on Emerging Topics in Computational Intelligence* 3, no. 6 (2018): 416-426.

[11] Nguyen, Van Nhan, Robert Jenssen, and Davide Roverso. "LS-Net: Fast single-shot line-segment detector." *Machine Vision and Applications* 32 (2021).

[12] Lee, Sang Jun, Jong Pil Yun, Hyeyeon Choi, Wookyong Kwon, Gyogwon Koo, and Sang Woo Kim. "Weakly supervised learning with convolutional neural networks for power line localization." In *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1-8. IEEE, 2017.

[13] Li, Bo, Cheng Chen, Shiwen Dong, and Junfeng Qiao. "Transmission line detection in aerial images: An instance segmentation approach based on multitask neural networks." *Signal Processing: Image Communication* 96 (2021): 116278.

[14] Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431-3440. 2015.

[15] Alwajih, Fakhraddin, Eman Badr, and Sherif Abdou. "Transformer-based Models for Arabic Online Handwriting Recognition." *International Journal of Advanced Computer Science and Applications* 13, no. 5 (2022).

[16] Gutiérrez Choque, Anyelo Carlos, Vivian Medina Mamani, Eveling Castro Gutiérrez, Rosa Núñez Pacheco, and José Ignacio Aguaded. "Transformer based Model for Coherence Evaluation of Scientific Abstracts: Second Fine-tuned BERT." (2022).

[17] Graham, Benjamin, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. "Levit: a vision transformer in convnet's clothing for faster inference." In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 12259-12269. 2021.

[18] Xiao, Tete, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. "Early convolutions help transformers see better." *Advances in neural information processing systems* 34 (2021): 30392-30400.

[19] Srinivas, Aravind, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. "Bottleneck transformers for visual recognition." In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16519-16529. 2021.

[20] d'Ascoli, Stéphane, Hugo Touvron, Matthew L. Leavitt, Ari S. Morcos, Giulio Biroli, and Levent Sagun. "Convit: Improving vision transformers with soft convolutional inductive biases." In *International Conference on Machine Learning*, pp. 2286-2296. PMLR, 2021.

[21] Guo, Jianyuan, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. "Cmt: Convolutional neural networks meet vision transformers." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12175-12185. 2022.

[22] Wang, Wenhai, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions." In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 568-578. 2021.

[23] Wang, Wenhai, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. "Pvt v2: Improved baselines with pyramid vision transformer." *Computational Visual Media* 8, no. 3 (2022): 415-424.

[24] Pan, Zizheng, Bohan Zhuang, Haoyu He, Jing Liu, and Jianfei Cai. "Less is more: Pay less attention in vision transformers." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 2, pp. 2035-2043. 2022.

[25] Pan, Zizheng, Jianfei Cai, and Bohan Zhuang. "Fast vision transformers with hilo attention." *Advances in Neural Information Processing Systems* 35 (2022): 14541-14554.

[26] Liu, Ze, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. "Swin transformer: Hierarchical vision transformer using shifted windows." In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012-10022. 2021.

[27] Ronneberger, Olaf, Philipp Fischer, and Thomas Brox. "U-net: Convolutional networks for biomedical image segmentation." In *Medical Image Computing and Computer-Assisted Intervention– MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pp. 234-241. Springer International Publishing, 2015.

[28] Hu, Jie, Li Shen, and Gang Sun. "Squeeze-and-excitation networks." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132-7141. 2018.

[29] Abdelfattah, Rabab, Xiaofeng Wang, and Song Wang. "Ttpla: An aerial-image dataset for detection and segmentation of transmission towers and power lines." In *Proceedings of the Asian Conference on Computer Vision*. 2020.

[30] Loshchilov, Ilya, and Frank Hutter. "Decoupled weight decay regularization." *arXiv preprint arXiv:1711.05101* (2017).

[31] Chen, Liang-Chieh, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. "Encoder-decoder with atrous separable convolution for semantic image segmentation." In *Proceedings of the European conference on computer vision (ECCV)*, pp. 801-818. 2018.

[32] Zhou, Zongwei, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. "Unet++: A nested u-net architecture for medical image segmentation." In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4*, pp. 3-11. Springer International Publishing, 2018.

[33] Jaffari, Rabeea, Manzoor Ahmed Hashmani, and Constantino Carlos Reyes-Aldasoro. "A novel focal phi loss for power line segmentation with auxiliary classifier U-Net." *Sensors* 21, no. 8 (2021): 2803.

[34] Zhou, Yichao, Haozhi Qi, and Yi Ma. "End-to-end wireframe parsing." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 962-971. 2019.

[35] Xue, Nan, Tianfu Wu, Song Bai, Fudong Wang, Gui-Song Xia, Liangpei Zhang, and Philip HS Torr. "Holistically-attracted wireframe parsing." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2788-2797. 2020.

# Dynamic Object Detection Revolution: Deep Learning with Attention, Semantic Understanding, and Instance Segmentation for Real-World Precision

Karimunnisa.shaik[1], Dr. Dyuti Banerjee[2], Dr. R. Sabin Begum[3], Narne Srikanth[4], Jonnadula Narasimharao[5], Prof. Ts. Dr. Yousef A.Baker El-Ebiary[6], Dr.E.Thenmozhi[7]

Assistant Professor, Department of Information Technology, Marri Laxman Reddy Institute of Technology and Management, Dundigal, Hyderabad, India-500043[1]
Assistant Professor, Department of Artificial Intelligence & Data Science (AI&DS), Koneru Lakshmaiah Education Foundation, Green Fields, Vaddeswaram, Guntur District, Andhra Pradesh, India-522302[2]
Assistant Professor, Department of Computer Applications, B. S. Abdur Rahman Crescent Institute of Science and Technology, Chennai, Tamil Nadu, India[3]
Assistant Professor, Department of CSE (AI & ML), RVR & JC College of Engineering Andhra Pradesh, India[4]
Associate Professor, Department of Computer Science and Engineering, CMR Technical Campus, Medchal, Hyderabad, Telangana, India – 501401[5]
Faculty of Informatics and Computing, UniSZA University, Malaysia[6]
Associate Professor, Department of Information Technology, Panimalar Engineering College, Chennai, India[7]

*Abstract*—Semantic and instance segmentation are critical goals that span a wide range of applications, from autonomous driving to object recognition in different fields. The existing approaches have limitations, especially when it comes to the difficult task of identifying and detecting minute things in intricate real-world situations. This work presents a novel method that uses a hybrid deep learning architecture with the Python programming language to smoothly combine semantic and instance segmentation. The suggested approach takes care of the pressing necessity in challenging real-world settings for accurate localization and fine-grained object detection. By combining the strengths of a Convolutional Neural Network (CNN) with a Bidirectional Long Short-Term Memory Network (BiLSTM), the hybrid model effectively achieves semantic segmentation by using sequential input and spatial information. A parallel attention method is smoothly included into the segmentation process to further improve the model's capabilities and enable the recognition of important object attributes. This study highlights the difficulties caused by changing environmental elements, highlighting the need for precise object location and understanding in addition to the complexities of fine-grained object detection. The suggested approach has an outstanding accuracy rate of 99.66%, outperforming existing approaches by 25.22%. This significant increase highlights the benefits that the hybrid design has over individual techniques and shows how effective it is at resolving issues that arise in dynamic real-world circumstances. The research highlights the importance of attention processes in deep learning and demonstrates how they might improve the specificity and accuracy of object detection and localization in intricate real-world scenarios. The improved performance of the suggested methodology is with well-known techniques like RCNN, CNN, and DNN, reaffirming its status as a reliable means of developing object localization and recognition in difficult situations.

*Keywords*—*Semantic segmentation; instance segmentation; convolutional neural network; bidirectional long short-term memory; attention mechanism*

## I. INTRODUCTION

Fine-grained object identification in visual computing is an important topic for addressing real-world issues by concentrating on minute distinctions and subtleties among comparable things, prioritizing precision in technologies such as manufacturing, self-driving vehicles, healthcare imaging, and wildlife preservation [1]. The phrase "challenging real-world environments" in this sense refers to a wide range of elements, such as obstructions, various lighting situations, various views, and abundantly generated scenery. Gaining an improved understanding of the minute differences that set items belonging to the identical grouping apart is the main objective of fine-grained object recognition. For example, fine-grained recognition of objects goes beyond standard object recognition to distinguish between certain dog kinds or cat species. standard object recognition could distinguish among a dog and a cat [2]. Innovative strategies are required to uncover small features and patterns that are invisible to the human eye, which frequently combine learning algorithms, deep neural networks, and characteristic engineering to achieve efficient analysis. Real-world scenarios with continually changing and unexpected characteristics provide considerable challenges to a seamless object recognition infrastructure. obstacles, visual obstacles generated by extra objects or features, and variations in light, such as shadows and views, can all complicate the identification process [3]. Recognition algorithms must be adaptable to variations in size, movement, and viewpoint caused by real-world views and complicated backgrounds. Sophisticated algorithms are required to discriminate between useful items and insignificant visual characteristics in busy and patterned

surroundings, providing smooth object detection in difficult conditions [4]. A comprehensive strategy comprising innovative algorithm creation, different training data sets, and in-depth knowledge of individual applications is required to increase the robustness and dependability of smooth object recognition algorithms for real-world applications.

Robots and AI systems use instance segmentation and semantic segmentation as core computer vision operations, allowing them to analyze visual input in a variety of circumstances and translate pixel-level data into useful information [5]. Semantic segmentation is an important approach for scene evaluation and object recognition that assigns pixel-level classification to certain sections of a picture. This allows computer systems to recognize object boundaries and grasp how the world is organized. Semantic segmentation improves its efficiency by distinguishing between object categories and individual occurrences, making it the cornerstone of scene interpretation [6]. In picture segmentation, each structure has a precise object mask, which is critical for robotics and AI systems to reliably discriminate and recognize instances of the same item category. This level of granularity is critical in complicated applications such as self-driving cars and healthcare imaging, where detecting unique instances in crowded settings is important. It is difficult to distinguish delicate elements from clutter limits segmentation algorithms in real-world scenarios. Conventional approaches fail to detect small-scale items accurately, while real-world difficulties like as obstructions, accessibility limits, and changing illumination conditions hamper object identification and localization even more. These issues demand unique tactics that go beyond existing approaches to provide dependable and seamless object assessment in difficult real-world scenarios. [7].

This paper introduces a hybrid deep learning architecture that combines semantic and instance segmentation capabilities to improve real-world segmentation algorithms. The framework employs attention processes to outperform existing strategies in tasks demanding precise object placement and identification. It integrates semantics and instance-level segmentation into a single framework, fixing specific flaws while maximizing benefits. The method employs numerous decoding paths and shared encoders to extract important information from input images, allowing the system to comprehend the image's overall meaning and design. This fundamental understanding facilitates scene design and relationships [8]. The framework employs a shared encoder to decrease processing needs, increase efficiency, and ensure a thorough grasp of input data. Strategic decoding branching switches the encoder, resulting in more accurate and consistent segmentation results while optimizing resource utilization and data transmission. The hybrid architecture uses attention processes to improve system efficiency, allowing it to concentrate on critical data locations, especially in complicated and crowded situations where fine-grained features are likely to be covered [9]. The system employs attention methods for segmentation, ensuring that computer resources are used efficiently and that meaningful instances

are distinguished from background noise. The combination of attention processes and the shared encoder increases flexibility by handling occlusions, angles of view, and illumination conditions. This integrated technique encourages a more detailed comprehension of settings, resulting in considerable improvements in accurate item location and identification in demanding real-world scenarios [10]. The hybrid deep learning framework enhances semantics and instance segmentation through the use of twin decoding approaches. It separates pictures into relevant segments for semantic segmentation by generating pixel-wise categorization projections, allowing for comprehensive image interpretation by identifying significant portions of the image [11]. The second decoder is particularly built for instance segmentation, resulting in exact object masking to discriminate between distinct instances of objects. This method improves higher-level semantic comprehension and enables fluid localization, resulting in a more complete and accurate contextual understanding.

The framework incorporates attention processes into decoding operations to improve its efficacy in demanding settings. These tactics guarantee that just the most significant parts of a picture are analysed, allowing the algorithm to focus on areas critical for object localization and recognition. This enhances its capacity to handle complicated things, respond fast to visual cues, and function effectively in tough conditions [12]. The hybrid architecture's performance was evaluated using a variety of real-world datasets, such as interior settings, animal photos, and metropolitan landscapes. The system outperformed advanced approaches in semantics and instance segmentation, proving its capacity to handle accurate item localization tasks in demanding environments such as interior settings and metropolitan landscapes [13]. The hybrid architecture, as a flexible solution, exhibits its capacity to excel across a number of domains, providing potential paths for improvements in fine-grained visual processing and supporting applications ranging from autonomous cars to medical imaging.

The key contributions of the article is,

- The study introduces a novel pre-processing step using Gaussian functions, enhancing feature extraction and smoothing in input data, contributing to improved overall model performance.

- The research leverages a fusion of Convolutional Neural Networks (CNN) and Bidirectional Long Short-Term Memory networks (BiLSTM) for semantic segmentation, enabling the model to capture both spatial and temporal dependencies in the data, resulting in more accurate and context-aware segmentation.

- The incorporation of attention mechanisms in instance segmentation significantly refines object delineation. The attention mechanism ensures that the model focuses on relevant regions, enhancing the precision and efficiency of instance segmentation in complex scenes.

- The study offers a nuanced view of the model's capabilities by introducing a thorough performance evaluation technique that takes into account variables including segmentation accuracy, computing efficiency, and resilience.

This article's remainder is organized as follows: In Section II, a summary of related research is provided. Section III presents the problem statement. The suggested approach's methodology and architecture are explained in Section IV of the article. The findings and subsequent discussion are covered in Section V. The conclusion is covered in Section VI.

## II.    RELATED WORKS

Automatic recognition of objects in 3D spaces is essential to working zone security, including ensuring adherence to building codes and averting accidents and fatalities at workplaces [14]. However, a number of difficulties, including correct three-dimensional object comprehension because of size changes and a lack of indicators in the three-dimensional environment, outstanding identification, exceptional segmentation of instances, and insufficient technical object databases with masking present significant challenges. These difficulties affect conventional manual techniques. The main discovery is to calculate pseudo-light recognition and reaching point clouds for three-dimensional object recognition using two-dimensional recognition of objects, segmentation of instances, and cameras perception. On the contrary hand, an upgraded cascading masks R-CNN is used to identify boundaries and masking for every two-dimensional object, while an enhanced characteristic pyramids system is presented for obtaining additional smooth object characteristics. An additional object classes with the boundaries and masking is introduced, and the AIM database for massive machinery identification is expanded. On the contrary hand, using deep learning, autonomous camera parameters estimation, a vision-based approach, and a spatial filtering, it is possible to retrieve pseudo-LiDAR point's clouds of objects generated by boundary boxes and masking from a monochromatic image. Numerous tests and evaluations reveal that the recently developed approach can recognize three-dimensional items and autonomously assess the security working areas. On the AIM information set, the suggested object recognition system produced the most advanced outcomes, and for the enhanced information set. The fresh framework will act as a starting point for studies regarding three-dimensional object recognition for additional three-dimensional positions.

The study extends Trans10K-v1, the initial significant database for translucent item differentiation, by providing an additional, smooth database known as Trans10K-v2 [15]. The novel dataset has a number of enticing advantages over Trans10K-v1, which only contains two constrained classes. (1) It is better suited for practical use since it comprises eleven smooth classes of translucent items that are frequently seen in the average household surroundings. (2) Compared to its predecessor, Trans10K-v2 presents additional difficulties for currently sophisticated segmentation techniques. Additionally, the Trans2Seg pipelines, unique transformer-based segmented pipelines, are suggested. Initially Trans2Seg's transformers encoders, which offers an overall responsive field as opposed to CNN's localized one, outperforms standard CNN systems in terms of performance. Subsequently, the study builds a collection of accessible designs as the question parameters of Trans2Seg's transformers decoder, where every instance acquires the statistical information of a particular group in the entire data set by treating semantic segmentation as an issue of dictionaries. Researchers compare Trans2Seg with a variety of than twenty current semantically segmented techniques, revealing that it greatly exceeds all CNN-based techniques and potentially solving the problem of translucent object segmentation.

For the creation of a successful computerized diagnostic framework, the identification and fragmentation of the new coronavirus infection of 2019 abnormalities using CT images are extremely important [16]. One of the finest options for creating such an instance is deep learning. However, a number of issues, involving as information variation, a wide range in the dimensions and form of the inflammation, lesions imbalances and an absence of annotations, restrict the effectiveness of DL techniques. In order to overcome these difficulties, a unique multitasking regression networks for categorizing COVID-19 lesions is put forward in the present research. The model's designation is MT-nCov-Net. The lesions identification is formulated as a multitasking structure extrapolation issue, allowing for the sharing of low-, medium-, and excellent-quality information across multiple assignments. In order to effectively acquire tiny and substantial lesion characteristics while minimizing the semantic disparity between various scale visualizations, a multiscale characteristic learning modules is introduced. This component captures the multiscale knowledge of semantics. Also included is a smooth lesions identification module that employs an adaptable dual-attention technique to find infected regions. The inflammatory regeneration modules then segment the infected regions using the obtained position mapping and the merged multiscale depictions. By reducing the COVID-19 area's form, MT-nCov-Net can properly fragment the inflammation through absorbing all of its features. MT-nCov-Net is empirically assessed on two open multisource information sets, and the results support its general efficacy above the state-of-the-art methodologies and show how well it works to solve the issues with COVID-19 diagnostics.

Due to the intricacy of plant images, segmenting plants is a difficult automated vision problem [17]. The study has to do increasingly more challenging activities in order to tackle various real-world issues. Instead of looking at the entire plant, the study must make distinctions between plant sections. The lack of information with thorough annotations is the main obstacle to multi-part segmentation. Actively annotating databases at the object component levels takes a lot of effort and money. The study suggests using pseudo-annotation with inadequately supervised training. In the article, researchers examine the minimally supervised training techniques currently in use and offer a productive pipeline for agrarian applications. It is made to deal with close object overlaps. For the plant component example and the entire plant scenario, the pipeline outperforms the starting point solutions by twenty-three percent and forty percent, respectively. To improve

simulation effectiveness, researchers also use instance-level enhancement. The goal of the method is to create a weakened segmentation masking that can be used to trim items from the source images and paste them onto fresh backdrops while the participant is being trained. On object component segmentation operations, the strategy gives us a fifty-five percent mAP gain over the initial state, and on entire plant segmentation operations, a seventy-two percent gain.

Recent developments in navigational autonomy have shown a greater preference for computation vision over conventional methods [18]. The majority of the places are built with individual movement in mind, which is how this works. They are therefore chocking full of visual indicators. In this way, the capacity to recognize objects visually is crucial for self-driving cars to prevent impediments while interacting with the outside environment. It is laborious and costly to gather information employing unmanned aerial vehicles that are capable of operating in the real environment. A database consisting of areas and conversations constitutes one of IT businesses' greatest resources as a result. Adopting an image-realistic three-dimensional simulation as the source of information is one way to address this issue. It is feasible to obtain a substantial quantity of information with this asset. Therefore, utilizing images from a frontend UAV camera moving through a three-dimensional simulator, the present study builds a collection of images for example segmentation. The Mask-RCNN, a cutting-edge deep learning approach, is used in the present research. The framework estimates per-pixel segmentation of instances from an image input. According to empirical findings, Mask RCNN performs better in the data we provide when improving the algorithm generated from the COCO dataset. Additionally, the intriguing findings in real-life information provide the suggested technique a solid generalization potential.

The inconsistent and fragmented nature of points cloud information in a non-Euclidean environment makes it difficult to fully employ smooth semantic properties [19]. A max pooling procedure is frequently employed to draw attention to particularly significant characteristics in the immediate area in attempting to illustrate the local characteristic for every centering point that is beneficial for improved contextual understanding. The max pooling method, however, ignores any additional geometrical local connections between every central location and its related neighborhood. In order to do this, the focused attention method shows promise in preserving node representations on graph-based information by paying consideration to every node in its immediate vicinity. Using stacking MLP units and a unique neural network called GA Point Net, the study offers a new method for analyzing point clouds. GA Point Net can acquire localized geometrical descriptions. To effectively utilize local characteristics, the study emphasizes various focus weights on every focal point's neighborhood. To completely retrieve localized geometrical patterns and improve the system's resilience, the study additionally mixes attentive characteristics with the regional identity characteristics produced by the focus grouping. The suggested GA Point Net structure performs at the cutting edge across the form

categorization and segmentation operations when evaluated on a variety of datasets used as benchmarks.

To be more precise, the characteristics with multiple scales are combined top-down to blend specifics and broad semantics to improve tiny item detection [20]. A pyramidal tiered attention mechanism made up of channels and spatial focus is created throughout the fusion of multi-scale information in order to see the item. Additionally, enhancing self-paced learning is used to direct the model to study challenging data. Two real-world datasets, an AMMW dataset and a publicly accessible PMMW dataset, are used for validating the technique that is suggested. The results of experiments show that the suggested strategy is preferable due to its versatility in identifying various devices concealed inside clothing while remaining harmless; the AMMW scanner has become a common tool for checking the safety of people in public settings in recent years. Nevertheless, due to intrinsic image noise, unknown item kind, and ambiguous status, it is very difficult to identify all concealed objects autonomously and precisely. Recent improvements in concealed object identification have been made by various current algorithms, particularly those that utilize deep learning. These techniques are effective for finding a few specific types of huge things, but they are ineffective for finding dim or imperfect hard objects. The state-of-the-art approaches are provided in this paper as a hidden finding of objects model with SPFAFN to handle this problem. SPFAFN obtains improved results on the two datasets with Maximum Precision.

The importance of object detection for systems that drive autonomously is rising [21]. Nevertheless, the use of existing object detectors to autonomously drive is constrained by their low accuracy or low inference ability. By integrating dilated convolutions and a SAM into the YOLOv3 design, a quick and precise object detector known as SA-YOLOv3 is suggested in this study. In order to enhance the accuracy of detection, the loss functional determined by GIoU and focused loss is rebuilt. With immediate inference, the suggested SA-YOLOv3 enhances YOLOv3 by 2.58 mAP and 2.63 mAP on the KITTI and BDD100K standards, respectively. Its improved compromise in terms of quickness and - accuracy in comparison to other cutting-edge detectors suggests that it is appropriate for use in self-driving vehicle applications. It is believed that this strategy, which integrates YOLOv3 with an attention mechanism for the first time will serve as a model for upcoming studies on autonomous vehicles.

To differentiate objectives from inferior categorization is the goal of smooth visual categorization [22]. It is thought to be a very challenging assignment since smooth images naturally exhibit significant inter-class variations and tiny intra-class variability. The majority of current methods employ CNN-based systems as characteristic extractors, which results in the generated exclusive areas including the majority of the object's components and missing to identify the truly crucial components. The perception transformers, which employ a mechanism for focus to gather broad context-relevant data to create a distant reliance on the desired object and then extracts more potent characteristics, has subsequently proven its effectiveness on a variety of imagine activities. The ViT approach might function inadequately in the

categorization of smooth images because it continues to place greater emphasis on overall coarse-textured data than localized smooth data. The study enhances the ViT framework and develops a concentration accumulating transformers to more effectively catch small variations across images. To improve communication between every transformer layering, the study specifically suggests an essential consideration accumulator. Additionally, the study provides an original data entropy selection to direct the algorithm in accurately obtaining discriminatory portions of the image. Numerous tests demonstrate that the suggested model architecture is capable of operating at an innovative modern facilities level on a number of widely used databases.

It's crucial and difficult to recognize facial expressions employing a DCNN [23]. Although significant efforts have been performed to improve FER accuracy using DCNN, earlier experiments have not yet been adequately generalized for use in practical situations. Conventional FER research is mostly restricted to regulated lab-posed fronted facing images, which do not encounter the difficulties associated with motion disintegrate head causes, obstructions, face distortions, and illumination under unregulated circumstances. The study suggested a SqueezExpNet architecture for extremely efficient FER systems that can tolerate fluctuations in the environment. It can benefit from global as well as local face data. The network was split into two distinct phases: a geometric concentration phase, which uses a SqueezeNet-like structure to collect localized highlighting data, and a geographic texturing phase, which consists of numerous compressed and extended levels to take use of the highest-level broad characteristics. To highlight significant localized facial areas, the study specifically developed the weighted masking of three-dimensional facial characteristics and employed element-wise combination with a geographical characteristic in the initial step. The subsequent phase of the network is then fed with the facial geographical imagine and its enhancements. To help overcome the unpredictability, a network of recurrent neural networks was created to cooperate with the emphasized data from two phases instead of only employing the SoftMax function, much like a classification system. The three top expressions databases were used in investigations encompassing simple and complex FER objectives. The technique produced cutting-edge findings and surpassed the current DCNN algorithms. Instantaneous FER could uncover possible applications in monitoring, well-being, and feedback mechanisms according to the created construction, chosen investigation approach, and published outcomes.

The thorough literature review explores a wide range of computer vision and deep learning subjects. By balancing camera perception, instance segmentation, and two-dimensional object recognition, the first study presents a novel hybrid deep learning architecture that adeptly tackles real-world complexities in three-dimensional object recognition. The results are promising and have significant implications for the identification of machinery and the development of autonomous driving technologies. Moreover, it presents Trans2Seg, a novel transformer-based segmentation pipeline that significantly outperforms the state-of-the-art CNN-based

techniques. To enhance and broaden the research landscape in this dynamic and ever-evolving domain, the study complements these breakthroughs by creating Trans10K-v2, a novel dataset featuring a variety of translucent object classes.

## III. PROBLEM STATEMENT

Deep convolutional neural networks (DCNNs) must operate accurately in uncontrolled, real-world contexts, making face emotion recognition an important yet difficult challenge. Advances in facial expression recognition (FER) accuracy have mostly concentrated on controlled lab environments, ignoring problems like head obstacles, motion disintegration, face distortions, and changing lighting. Effective generalization in such real-world scenarios is not possible with the current DCNN-based models [23]. In order to overcome these shortcomings, the SqueezExpNet architecture is presented in this paper. It is intended for very effective FER systems that can adapt to changes in the environment and make use of both local and global face data. The network is divided into two phases: a geographic texturing phase and a geometric concentration phase. For enhanced performance, a network of recurrent neural networks and weighted masking of three-dimensional face characteristics are included. The suggested technology outperforms existing DCNN algorithms and has potential uses in feedback systems, monitoring, and wellbeing applications.

## IV. PROPOSED HYBRID LEARNING ARCHITECTURE WITH ATTENTION MECHANISMS

The approach entails putting forth a unique architecture that combines instance and semantic segmentation to improve precise item localization and recognition in difficult real-world circumstances. The integration recognizes that although semantic segmentation delivers high-level contextual comprehension, instance segmentation offers detailed instance-level information. The CNN-BiLSTM for semantic segmentation is a novel addition that makes use of the spatial context and sequence data gathered by BiLSTMs and CNNs. In particular in complicated situations, the addition of attention processes, such as segmentation, assists in focusing on key visual areas for exact instance classification. In addition to demonstrating the architecture's superiority over stand-alone approaches and emphasizing the importance of attention processes in boosting instance segmentation's accuracy and reliability, the research makes sure to provide a thorough review of the architecture through quantitative and qualitative assessments. It is depicted in Fig. 1.



Fig. 1. Proposed methodology.

## A. Data Collection

The VEDAI dataset, which supports the growth of fine-grained vehicle recognition algorithms in remote sensing images, is a dataset for fine-grained vehicle identification. The total number of images and occurrences in VEDAI is 1210, and each image has a 1024 by 1024 pixel resolution. The order of occurrences in this collection can be considered sparse, as can be observed by the quantity of instances and images [25].

## B. Pre-Processing using Gaussian Function

When preprocessing an image using a Gaussian function, a Gaussian filter is applied to it to assist smooth it out and remove noise. The kernel or mask of the filtering procedure is the bell-shaped Gaussian function, a mathematical function. The filter operates using a Gaussian kernel, and each pixel is assigned a weight based on how near to its neighbors it is. The weights are computed using the values of the Gaussian function at each location along the kernel. By changing the Gaussian filter's parameters, such as the kernel size and standard deviation, the amount of image smoothing may be changed. While bigger kernel sizes and greater standard deviations provide more comprehensive smoothing, smaller kernel sizes and lower standard deviations maintain finer features. The final preprocessed image has a smoother overall look, less noise, and fewer high-frequency features. Gaussian filtering, a preprocessing technique frequently used in image processing, is particularly useful for tasks like demising, feature extraction, and increasing the quality of images. Because the Gaussian filter can blur and smooth images well, it is a crucial tool in computer vision and image processing applications. This filter, which gets its name from the Gaussian function, is frequently used to improve images and reduce noise. Through the use of a Gaussian kernel to convolve the input picture, the filter reduces high-frequency noise while maintaining important image characteristics. Its natural smoothing ability helps to lessen pixel-level differences, producing an output that is more aesthetically pleasant and visually cohesive. The blurring effect of the Gaussian filter also helps to emphasize important structures while minimizing unimportant details, which makes it useful in applications like edge detection and feature extraction. Because of its adaptability, the filter may be used in a wide range of domains, such as object identification, robotics, and medical imaging. As such, it is an essential tool for preprocessing and picture refinement in a variety of settings. The Gaussian function equation is shown below in Eq. (1).

$$H(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{x^2}{2\pi\sigma^2}} \qquad (1)$$

where the standard deviation of the distribution is denoted as $\sigma$. The distribution is assumed to have a mean of 0.

## C. Hybrid CNN-BiLSTM for Semantic Segmentation

Considering it effectively handles image interpretation issues by combining the capabilities of CNN and BiLSTM, the Hybrid CNN-BiLSTM architecture is essential to semantic segmentation. Through the combination of the sequential information-taking capacity of the BiLSTM and the spatial feature-capturing capability of the CNN, this innovative approach allows the model to interpret both contextual and spatial subtleties in the image data. These two components work together in the Hybrid CNN-BiLSTM architecture to significantly increase semantic segmentation accuracy and precision, particularly where precise localization and fine-grained object recognition are crucial. The resilience and reliability of semantic segmentation tasks are enhanced by this design, which makes it possible for the model to function effectively in dynamic and complicated real-world environments. This makes it vital for a wide range of applications, including autonomous driving, computer vision, and medical imaging.

A convolutional neural network (CNN) is used in this step of feature extraction and is a crucial step in deep learning for pattern and recognition of image applications. CNN is developed in order to autonomously determine organizational and discriminating qualities from the original input images. In a series of convolutional layers, tiny filters are organized over the input image to collect regional trends and attributes. These characteristics capture the edges, the surfaces, and shapes which make up the visual environment. Then, the feature maps are down sampled by pooling layers, which lessens the computational cost while maintaining important data. Fully connected layers accept the results of the learnt attributes and put them into effect future classification or similar activities.

CNNs have shown outstanding results in a range of applications involving computer vision, demonstrating modern performance in problems including object detection, image segmentation, and analysis of images in healthcare. This is because of their ability to voluntarily pick up and remember important aspects from visuals. By merging information from different methods, the CNN technique maximizes the benefits of each modality while compensating for its limitations and providing a more complete knowledge. Recurrent neural networks of the Bidirectional Long Short-Term Memory (BiLSTM) variety handle sequential information while simultaneously considering into account both past as well as future circumstances. It has both forward and backward LSTM components that create hidden states and collect data from components that come before and after them in the sequence. When combined with convolutional neural networks (CNNs), BiLSTMs can improve pixel-wise image categorization in applications like semantic segmentation. Bidirectional connections are utilized into consideration by BiLSTMs to enhance the framework's comprehension of intricate spatial and contextual interactions, producing more precise and context-sensitive semantic segmentation that benefits applications related to computer vision. The diagrammatic representation of the proposed hybrid CNN-BiLSTM method is depicted in Fig. 2 and the equations of CNN-BiLSTM are depicted in Eq. (2) to Eq. (7) below:

Fig. 2. CNN-BiLSTM framework.

$$g_v = \sigma(N_g y_d + K_g g_{v-1} + d_g) \qquad (2)$$

$$h_v = tan_g(N_h y_v + K_h g_{v-1} + d_h) \qquad (3)$$

$$j_v = \sigma(N_j y_v + K_q g_{v-1} + d_j) \qquad (4)$$

$$p_v = \sigma(N_i y_v + K_p g_{v-1} + d_p) \qquad (5)$$

The current qv state can be calculated by (6).

$$q_v = g_v \times q_{v-1} + j_v \times h_v \qquad (6)$$

$$x_v = g_v = p_v \times tan_g(q_v) \qquad (7)$$

Here, Kg, Kh, Kq, Kp signifies the weight matrices of the previous short-term state gv-1. Ng, Nh, Nj, Ni signifies the weight matrices of the current input state yd, dh, dg, dj, and dp are labelled as the bias terms, qv-1 characterizes the preceding long-term state.

### D. Attention Mechanism for Instance Segmentation

The Attention Mechanism is a fundamental component of instance segmentation and serves as a basis for enhancing and optimizing object recognition in intricate visual datasets. The primary function of this method is to enhance the accuracy and precision of instance segmentation by allowing the model to recognize and concentrate on certain object features and regions. Even in congested or densely populated environments, the model can efficiently separate and isolate individual object instances because to its dynamic capacity to balance the value of different components in the visual input. This adaptable approach proves its worth in a multitude of domains, ranging from robotics and object recognition to the intricate realm of medical imaging, where it enhances instance segmentation accuracy and consistency. This paves the door for innovative developments in these domains by laying the foundation for a richer comprehension of complicated visual data and the things it includes.

The attention model improves the CNN by maintaining its context-relevant characteristics. In the prior-based model, each block's attributes are integrated into the ones from the layer below it. This method equally weights each attribute acquired from the previous CNN blocks. To learn precise feature values, important characteristics from the previous blocks must be given a high weight relative to other features. In order to facilitate the acquisition and selection of noteworthy qualities from previous blocks, a mechanism tracking attention was consequently introduced to the CNN architecture. This model generates an attention mask that equalizes the relative importance of spatial characteristics on that feature map. Leveraging an attention mechanism throughout blocks, the CNN architecture generates a weighted function for simulating activations from the prior blocks. The connections from the previous blocks that were skipped were then weighed throughout the depth axis for every single pixel in that layer's spatial range [26].

Two operational channels, H(x) and S'(x) are used to guide the layer of convolution to produce 'x' output in both the first and subsequent blocks. H(x) shows the set of methods that were implemented to take the input value "x" and just feed it forward to the block after it. The group of techniques collectively referred to as S'(x) has weighting with attention and skipping the 'x' through convolutional and maximum-pooling layers. The balanced summation is used to obtain the outputs G(x) from the CNN block and is shown in Eq. (8).

$$G(X) = H(X) + S'(X) \qquad (8)$$

Eq. (9) is used to determine the functional route S'(x).

$$S'(X) = S'(X) * \varphi \qquad (9)$$

where S(x)'s spatial dimensions and the attention weight matrix's parameters are equivalent. The appropriate cross-section of S(x) is multiplied point-wise (broadcast throughout the depths) by the attention matrix's weights "." CNN can incorporate the inputs from the current moment and its results from the previous instant to automatically allocate the weighting for each element of the network as a whole by including an attention method. Important image details may be focused on to improve the classification's precision and adaptability. Based on this, CNN's attention concept is added to form the Attention-CNN model. Pay attention carefully: To identify and classify the framework and airborne particulates in SEM images, CNN is used. The attention-CNN architecture consists of the input, a convolution, attention to detail, full connection, and output layers. The layer that receives input is made up of four nodes, which are the labelled images of four separate kinds of particles, as the input data is a pictorial representation of four distinct particle types. Each of the convolution layers of the four phases that together make up a single layer of convolution is followed by a layer of attention in order to successfully accomplish weight transportation. The dimensions of the convolution kernels are 8x3x3, 16x5, 32x3, and 32x3 for each layer, accordingly. A pooling layer connects the first and last convolution layers. The final result of the layer used for convolution serves as the layer's input, and the total number of nodes in the full connections layer is set to 64. The maximum number of the output layer nodes is 4, and the output layer classifies the smallest particles into four groups.

Fig. 3. Attention mechanism integration in CNN.

The CNN Attention Mechanism for Integration is shown in Fig. 3. A nonlinear relationship might be inserted amongst the different layers of an ensemble of neurons by changing the function that causes activation. The network's output no longer looks linear, which increases the network's expression and allows it to fit a wider range of patterns. The Attention-CNN model uses two activating coefficients, Relu (rectified linear unit) as well as, for its concealed and outputs components, respectively. Relu can deal with elevation dispersion throughout the element transfer process. When the value of the Relu function is greater than 0, its derivative is 1. It is simple to identify the gradient and may greatly accelerate the gradients' downward speed of convergence. Eq. (10) shows the Relu function's formulation.

$$ReLu = max(0, X) \tag{10}$$

By transferring the outputs of several neurons to the coordinates (0, 1), Softmax is able to classify data in a variety of ways. Assuming there is one, $j$ denote the last component of an input array; the softmax value assigned to that component is determined by

$$G_j = \frac{s^j}{\sum_{i=1}^{k} s^j} \tag{11}$$

where, k stands for all of the input items. Since both the first-order and second-order instant means of the gradient are completely used by the Adam optimization approach, the forward momentum component is taken into account during the updating procedure. Adam's computation is shown in Eq. (12) to Eq. (16):

$$u_t = \alpha_1 \delta_{t-1} + (1 - \alpha_1)k_t \tag{12}$$

$$n_t = \alpha_2 \delta_{t-1} + (1 - \alpha_2)k_t^2 \tag{13}$$

$$\hat{u}_t = \frac{u_t}{1 - \beta_1^t} \tag{14}$$

$$\hat{n}_t = \frac{n_t}{1 - \delta_1^t} \tag{15}$$

$$k_{t+1} = k_t - \frac{\partial}{\sqrt{\hat{n}_t} + \varepsilon} \hat{u}_t \tag{16}$$

where $u_t$ is the first-order moment estimates that nt is the second-order momentum term, $\alpha_1$, $\alpha_2$ are actually dynamic values, $k_t$ is the gradient of the cost operates after t iterations, $u_t$ is the first moment's correction value, $n_t$ is the second moment's correction value, $k_t$ is the model's variables, and is a small amount that can circumvent the zero denominators. In neural network [24] training, the loss of functions is used to quantify the difference between the predicted result and the actual value. In addition, the effectiveness of the computational framework is evaluated using this component as a benchmark. The cross-entropy cost functioning, which may be thought of as the loss of function for Attention-CNN in Eq. (17),

$$H = -\frac{1}{n} \sum_l l_j(\rho(a) - b) \tag{17}$$

where, y is the resultant value, b is the actual value, n is the total of the samples l is the sample, and n is the sample. The following formula is used to determine the gradient.

$$\frac{\partial b}{\partial a} = \frac{1}{n} \sum_l l_j(\omega(z) - b) \tag{18}$$

where the error between the output and the actual value is $\omega(z) - b$ [27].

## V. RESULTS AND DISCUSSIONS

To improve accurate item localization and recognition in challenging real-world scenarios, the approach proposes a novel architecture that combines instance and semantic segmentation. The integration acknowledges that while instance segmentation provides detailed instance-level information, semantic segmentation provides high-level contextual comprehension. A novel addition that leverages the spatial context and sequence data collected by CNNs and BiLSTMs is the CNN-BiLSTM for semantic segmentation. The addition of attention processes, like segmentation, helps to focus on important visual areas for precise instance classification, especially in complex scenarios. The research ensures that the architecture is thoroughly reviewed through quantitative and qualitative assessments, highlighting its superiority over stand-alone approaches and highlighting the role of attention processes in improving the accuracy and reliability of instance segmentation.

### A. Performance Metrics

Training and Testing Accuracy: An indicator of a machine learning model's success during the training stage is training accuracy. It displays the percentage of instances (or samples) in the training dataset that were properly predicted in relation to all of the occurrences in that dataset. In other words, the degree to which the model's predictions match the actual labels for the data it was trained on is indicated by the degree of training accuracy.



Fig. 4. Training and testing accuracy.

In this Fig. 4, accuracy refers to a performance parameter that assesses the percentage of fine-grained objects that the proposed hybrid deep learning architecture successfully recognized and localized in difficult real-world contexts. The model's excellent precision and efficacy in successfully recognizing and localizing items with complex features and difficult backgrounds are indicated by the accuracy of 99.66% that was attained.

Training and Testing Loss: A machine learning model's training loss, often referred to as the objective or cost function, is a quantifiable indicator of how well the model is doing. It shows the difference between the actual target values (ground truth) found in the training dataset and the projected values produced by the model. Making the model's predictions as near to the actual values as is practical is the main objective of training; this is to minimize this loss.



Fig. 5. Testing loss.

In Fig. 5, the gap between anticipated and ground truth values during training is measured as loss, which is shown graphically, and iteratively evolves over time for the hybrid deep learning architecture. As the model is trained, the loss curve shows how the architecture is coming together to minimize mistakes and enhance its capacity to precisely recognize and localize fine-grained objects in challenging real-world circumstances.

ROC Curve: In binary classification tasks, a graphical depiction known as the Receiver Operating Characteristic (ROC) curve is frequently used to assess how well a machine learning model is doing. As the discriminating threshold for distinguishing positive and negative occurrences is changed, it demonstrates the trade-off between the TPR, also known as sensitivity or recall, and the FPR.



Fig. 6. ROC curve.

In Fig. 6, the best threshold for detection assessments is present; the ROC assesses the model's capacity to discriminate between the presence and absence of objects. ROC values over a certain threshold indicate a more precise object identification model. The true positive rate, also referred to as sensitivity or recall, is shown on the vertical axis, while the false positive rate is represented on the horizontal axis. The proportion of actual positive events that the classification model correctly identified is measured by the true positive

rate, while the percentage of no object data that is incorrectly classified as an object is shown by the false positive rate. These rates for various classification thresholds, which specify the point at which the model classifies a data point as an event or a non-event, are plotted to create the ROC curve.

Matthews Correlation Co-efficient (MCC): The MCC is one of the most popular measures for classification effectiveness. It is generally recognized as a reliable estimate that may be used even when class sizes vary significantly. Eq. (19) contains the formula for the Matthews correlation coefficient.

$$MCC = \frac{T_P T_N - F_P F_N}{\sqrt{(T_P + F_P) - (T_P + F_N)(T_N + F_P)(T_N + F_N)}} \quad (19)$$

Negative Predictive Value (NPV): The subject-to-outcome ratio calculates the percentage of patients who have really poor test findings overall. The NPV measures the proportion of times each forecast was entirely wrong. Eq. (20) has the formula.

$$NPV = \frac{T_N}{T_N + F_N} \quad (20)$$

TABLE I. COMPARISON OF MCC AND NPV

| Methods | NPV (%) | MCC (%) |
|---|---|---|
| RCNN | 89.82 | 69.47 |
| CNN | 85.67 | 57.03 |
| DNN | 79.58 | 38.75 |
| **Hybrid DL-Attention Mechanism** | **93.77** | **73.12** |

Based on their NPV and MCC performance indicators, the various techniques are compared in Table I. The RCNN, CNN, DNN, and a hybrid deep learning (DL) model with an attention mechanism are the four techniques that are assessed. The comparison table's findings highlight the importance of the Hybrid DL-Attention Mechanism as a useful strategy for the assigned job, surpassing more established approaches like RCNN, CNN, and DNN in terms of NPV and MCC. The results highlight the rising significance of attention processes in deep learning by demonstrating their capacity to improve correlation metrics and prediction accuracy across a range of applications. It is shown in Fig. 7.



Fig. 7. Comparison of NPV and MCC.

False Positive Rate: The percentage of situations, where a favorable outcome was predicted but it did not materialize. It is illustrated in Eq. (21).

$$FPR = \frac{F_P}{T_N + F_P} \quad (21)$$

False Negative Rate: Eq. (22) displays the percentage of positive situations that were expected to be negative but ended up being positive.

$$FNR = \frac{F_N}{T_P + F_N} \quad (22)$$

TABLE II. COMPARISON OF FPR AND FNR

| Methods | FPR (%) | FNR (%) |
|---|---|---|
| RCNN | 10.17 | 20.35 |
| CNN | 14.32 | 28.64 |
| DNN | 20.41 | 40.82 |
| **Hybrid DL-Attention Mechanism** | **3.51** | **4.44** |

The contrast provided Table II highlights the FPR and FNR performance of the Hybrid DL-Attention Mechanism. The results highlight the potential of deep learning's attention mechanisms by showing how they may considerably improve the model's performance in tasks that call for striking a careful balance between reducing false alarms and missed detections. By demonstrating the advantages of attention processes in optimizing detection and classification models for practical applications, this research adds to the larger area of computer vision (see Fig. 8).



Fig. 8. Comparison of FPR and FNR.

Accuracy: Accuracy is used to evaluate the system model's performance as a whole. Every interaction can be properly foreseen is its central tenet. Eq. (23) provides the precision.

$$Accuracy = \frac{T_{Pos} + T_{Neg}}{T_{Pos} + T_{Neg} + F_{Pos} + F_{Neg}} \quad (23)$$

Precision: Precision also describes how closely two or more computations resemble one another in addition to being right. The link between precision and accuracy shows how often a judgment may be formed. Precision may be calculated using (24).

$$P = \frac{T_{Pos}}{T_{Pos} + F_{Pos}} \qquad (24)$$

Recall: The percentage of all pertinent results that were effectively sorted by the procedures is known as recall. The suitable positive for such numbers is determined by dividing the genuine positive by the erroneous negative values. It is referenced in Eq. (25).

$$R = \frac{T_{Pos}}{T_{Pos} + F_{Neg}} \qquad (25)$$

F1 Score: Accuracy and recall are combined in the F1-Score calculation. Eq. (26) computes the F1-Score using precision and recall.

$$F1 - score = \frac{2 \times precision \times recall}{precision + recall} \qquad (26)$$

A comparison of many object identification techniques, including RCNN, CNN, DNN, and a hybrid DL architecture

with attention mechanisms, is shown in Table III. Key performance indicators including accuracy, precision, recall, and F1-Score are used to assess these approaches. With an astounding accuracy of 99.66%, the hybrid DL architecture with Attention Mechanisms emerges as the most accurate technique. The accuracy, recall, and F1-Score values for this approach are also superior, coming in at 98.12%, 97.65%, and 96.54%, respectively. The outcomes highlight how combining attention mechanisms into a hybrid DL architecture has a substantial influence and leads in remarkable object detection performance. Overall, the evaluation's findings highlight the promise of cutting-edge strategies that make use of attention processes to improve object detection resilience and accuracy in difficult situations (see Fig. 9).

TABLE III. COMPARISON OF PERFORMANCE METRICS

| Methods | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| RCNN | 86.43 | 79.64 | 79.64 | 79.64 |
| CNN | 80.90 | 71.35 | 71.35 | 71.35 |
| DNN | 72.78 | 71.35 | 71.35 | 59.17 |
| Hybrid DL-AM | **99.66** | **98.12** | **97.65** | **96.54** |



Fig. 9. Comparison of performance metrics.

## VI. DISCUSSION

The paper addresses recent advances in computer vision, focusing on the application of deep learning algorithms for workplace safety, automated diagnostic frameworks for COVID-19 abnormalities [16], and plant segmentation. It also looks at navigational autonomy, proposing the use of 3D simulations in self-driving automobiles and Mask-RCNN for picture segmentation. The study also examines object detection and introduces SA-YOLOv3 for self-driving applications [21]. It also emphasizes the need of seamless visual classification through enhanced perception transformers. The study on facial expression recognition proposes the SqueezExpNet architecture [23], which efficiently solves real-world issues such as motion, obstacles,

and lighting changes. The study emphasizes the ongoing advancement of computer vision techniques, as well as their potential influence on safety, healthcare, agriculture, autonomous systems, and other fields.

## VII. CONCLUSION AND FUTURE WORK

This research examines the profound challenges associated with precise object localization and identification in complex real-world scenarios. The subject of computer vision has advanced significantly with the introduction of the new hybrid deep learning architecture. By integrating semantic and instance segmentation approaches with attention mechanisms to maximize their respective strengths, the strategy has demonstrated exceptional performance. Comprehensive tests and studies have clearly shown significant improvements over

conventional methods in terms of precision, memory, and overall correctness. Among the notable achievements are the integration of CNN and BiLSTM for semantic segmentation and the astute use of attention processes in instance segmentation. The model can now find fine-grained objects with the maximum accuracy, even in dynamic and complicated real-world scenarios, thanks to these advancements. This discovery has extensive implications for many disciplines, including medical imaging, robotics, autonomous driving, and more, that depend on reliable and consistent object localization. With the increasing complexity of real-world environments, the hybrid design in conjunction with attention mechanisms has opened up new and intriguing possibilities to enhance the adaptability and reliability of computer vision systems. Further study in this hybrid design in several domains and situations is a potential avenue. Particular difficulties arise with applications in industrial automation, agriculture, and underwater research. More substantial advancements in computer vision could result from evaluating the architecture's resilience and flexibility under these conditions. This forthcoming study aims to scrutinize the versatility of the hybrid technique, exploring its capacity to redefine and revolutionize object localization across diverse real-world scenarios.

REFERENCES

[1] Z. Yang, X. Yang, M. Li, and W. Li, "Small-sample learning with salient-region detection and center neighbor loss for insect recognition in real-world complex scenarios," Computers and Electronics in Agriculture, vol. 185, p. 106122, 2021.

[2] A. Abdelreheem, U. Upadhyay, I. Skorokhodov, R. Al Yahya, J. Chen, and M. Elhoseiny, "3dreftransformer: Fine-grained object identification in real-world scenes using natural language," in Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2022, pp. 3941–3950.

[3] S. Xiong, G. Tziafas, and H. Kasaei, "Enhancing Fine-Grained 3D Object Recognition using Hybrid Multi-Modal Vision Transformer-CNN Models," in IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2023), 2023.

[4] Y. Xi, W. Jia, Q. Miao, X. Liu, X. Fan, and H. Li, "FiFoNet: Fine-Grained Target Focusing Network for Object Detection in UAV Images," Remote Sensing, vol. 14, no. 16, p. 3919, 2022.

[5] Y. Chu et al., "A Fine-Grained Attention Model for High Accuracy Operational Robot Guidance," IEEE Internet of Things Journal, vol. 10, no. 2, pp. 1066–1081, 2022.

[6] S. Ye et al., "CDLT: A Dataset with Concept Drift and Long-Tailed Distribution for Fine-Grained Visual Categorization," arXiv preprint arXiv:2306.02346, 2023.

[7] J. Song, L. Miao, Q. Ming, Z. Zhou, and Y. Dong, "Fine-Grained Object Detection in Remote Sensing Images via Adaptive Label Assignment and Refined-Balanced Feature Pyramid Network," IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, vol. 16, pp. 71–82, 2022.

[8] T. A. Abir, E. Kuantama, R. Han, J. Dawes, R. Mildren, and P. Nguyen, "Towards Robust Lidar-based 3D Detection and Tracking of UAVs," in Proceedings of the Ninth Workshop on Micro Aerial Vehicle Networks, Systems, and Applications, 2023, pp. 1–7.

[9] L. Dodds, I. Perper, A. Eid, and F. Adib, "A Handheld Fine-Grained RFID Localization System with Complex-Controlled Polarization," arXiv preprint arXiv:2302.13501, 2023.

[10] D. Rathnayake, M. Radhakrishnan, I. Hwang, and A. Misra, "LILOC: Enabling precise 3D localization in dynamic indoor environments using LiDARs," 2023.

[11] M. Cimdins, S. O. Schmidt, F. John, M. Constapel, and H. Hellbrück, "MA-RTI: Design and Evaluation of a Real-World Multipath-Assisted

[12] I. Kar, S. Mukhopadhyay, and B. Guha, "A Dual Fine Grained Rotated Neural Network for Aerial Solar Panel Health Monitoring and Classification," in International Conference on Data Management, Analytics & Innovation, Springer, 2023, pp. 457–477.

[13] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn, "Learning fine-grained bimanual manipulation with low-cost hardware," arXiv preprint arXiv:2304.13705, 2023.

[14] "Deep learning-based object identification with instance segmentation and pseudo-LiDAR point cloud for work zone safety - Shen - 2021 - Computer-Aided Civil and Infrastructure Engineering - Wiley Online Library." https://onlinelibrary.wiley.com/doi/abs/10.1111/mice.12749 (accessed Aug. 03, 2023).

[15] E. Xie et al., "Segmenting Transparent Object in the Wild with Transformer." arXiv, Feb. 23, 2021. doi: 10.48550/arXiv.2101.08461.

[16] W. Ding, M. Abdel-Basset, H. Hawash, and O. M. ELkomy, "MT-nCov-Net: A Multitask Deep-Learning Framework for Efficient Diagnosis of COVID-19 Using Tomography Scans," IEEE Transactions on Cybernetics, vol. 53, no. 2, pp. 1285–1298, Feb. 2023, doi: 10.1109/TCYB.2021.3123173.

[17] S. Mukhamadiev, S. Nesteruk, S. Illarionova, and A. Somov, "Enabling Multi-Part Plant Segmentation with Instance-Level Augmentation Using Weak Annotations," Information, vol. 14, no. 7, Art. no. 7, Jul. 2023, doi: 10.3390/info14070380.

[18] F. X. Viana, G. M. Araujo, M. F. Pinto, J. Colares, and D. B. haddad, "Aerial Image Instance Segmentation Through Synthetic Data Using Deep Learning," Learn. Nonlin. Mod., vol. 18, no. 1, pp. 35–46, Sep. 2020, doi: 10.21528/lnlm-vol18-no1-art3.

[19] C. Chen, L. Z. Fragonara, and A. Tsourdos, "GAPointNet: Graph attention based point neural network for exploiting local feature of point cloud," Neurocomputing, vol. 438, pp. 122–132, May 2021, doi: 10.1016/j.neucom.2021.01.095.

[20] X. Wang et al., "Self-Paced Feature Attention Fusion Network for Concealed Object Detection in Millimeter-Wave Image," IEEE Transactions on Circuits and Systems for Video Technology, vol. 32, no. 1, pp. 224–239, Jan. 2022, doi: 10.1109/TCSVT.2021.3058246.

[21] D. Tian et al., "SA-YOLOv3: An Efficient and Accurate Object Detector Using Self-Attention Mechanism for Autonomous Driving," IEEE Transactions on Intelligent Transportation Systems, vol. 23, no. 5, pp. 4099–4110, May 2022, doi: 10.1109/TITS.2020.3041278.

[22] Q. Wang, J. Wang, H. Deng, X. Wu, Y. Wang, and G. Hao, "AA-trans: Core attention aggregating transformer with information entropy selector for fine-grained visual classification," Pattern Recognition, vol. 140, p. 109547, Aug. 2023, doi: 10.1016/j.patcog.2023.109547.

[23] A. R. Shahid and H. Yan, "SqueezExpNet: Dual-stage convolutional neural network for accurate facial expression recognition with attention mechanism," Knowledge-Based Systems, vol. 269, p. 110451, Jun. 2023, doi: 10.1016/j.knosys.2023.110451.

[24] P. Achlioptas, A. Abdelreheem, F. Xia, M. Elhoseiny, and L. Guibas, "Referit3d: Neural listeners for fine-grained 3d object identification in real-world scenes," in Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16, Springer, 2020, pp. 422–440.

[25] X. Sun et al., "FAIR1M: A Benchmark Dataset for Fine-grained Object Recognition in High-Resolution Remote Sensing Imagery." arXiv, Mar. 24, 2021. Accessed: Aug. 07, 2023. [Online]. Available: http://arxiv.org/abs/2103.05569.

[26] K. R., H. M., S. Anand, P. Mathikshara, A. Johnson, and M. R., "Attention embedded residual CNN for disease detection in tomato leaves," Applied Soft Computing, vol. 86, p. 105933, Jan. 2020, doi: 10.1016/j.asoc.2019.105933.

[27] C. Yin, X. Cheng, X. Liu, and M. Zhao, "Identification and Classification of Atmospheric Particles Based on SEM Images Using Convolutional Neural Network with Attention Mechanism," Complexity, vol. 2020, p. e9673724, Sep. 2020, doi: 10.1155/2020/9673724.

Device-Free Localization System," Sensors, vol. 23, no. 4, p. 2199, 2023.

# Improved Algorithm with YOLOv5s for Obstacle Detection of Rail Transit

Shuangyuan Li[1], Zhengwei Wang[2], Yanchang Lv[3], Xiangyang Liu[4]

Information Construction Office, Jilin Institute of Chemical Technology, Jilin, China[1,3]

School of Information and Control Engineering, Jilin Institute of Chemical Technology, Jilin, China[2,4]

*Abstract*—As an infrastructure for urban development, it is particularly important to ensure the safe operation of urban rail transit. Foreign object intrusion in urban rail transit area is one of the main causes of train accidents. To tackle the obstacle detection challenge in rail transit, this paper introduces the CS-YOLO urban rail foreign object intrusion detection model. It utilizes the improved YOLOv5s algorithm, incorporating an enhanced convolutional attention CBAM module to replace the original YOLOv5s algorithm's backbone network C3 module. In addition, the KM-Decoupled Head is proposed to decouple the detection head, and SIoU is applied as the loss function. Tested on the WZ dataset, the average accuracy increased from 0.844 to 0.893 .The research method in this paper provides a reference basis for urban rail transit safety detection, which has certain reference value.

*Keywords—Railroad track intrusion detection; CBAM (Convolutional Block Attention Module) attention; activation function; decoupling probe; loss function*

## I. INTRODUCTION

The safety environment of urban rail transit is crucial for ensuring the safety of rail transit operations and the well-being of passengers.

With the gradual improvement of safety requirements in the urban rail transportation industry, the number of safety accidents in rail transportation has significantly reduced. However, due to natural and human factors leading to foreign object intrusion on railroad tracks, these incidents exhibit sudden and unpredictable characteristics. Relying solely on manpower cannot ensure comprehensive and timely detection. Therefore, it is necessary to establish a corresponding railroad track foreign object intrusion detection system for real-time monitoring of the running lines. This system aims to provide accurate alarm information to duty personnel, ensuring the safety of train operations. The detection of foreign object intrusions on urban rails has been a focal point of research, and an effective intrusion detection method for urban rail is of great significance for ensuring the overall safety of urban rail operations.

At present, there are two main methods for urban rail transit intrusion detection: contact detection and non-contact detection. Contact detection requires a large amount of hardware as a support and is troublesome to install, which is not easy to use on a large scale, and at the same time, when the equipment detects the intrusion of foreign objects it cannot be disposed of in time, which will seriously affect the safety of urban transportation. However, computer vision is an efficient non-contact intrusion detection method that is widely used in different transportation industries. It has the advantages of easy maintenance and intuitive results. However, the complexity and variability of urban environments and disturbances such as bad weather can lead to false alarm problems. Fortunately, with the development of deep learning algorithms, it has been possible to achieve high detection accuracy and reduce the false alarm rate to a certain extent. Deep learning algorithms excel at automatically learning complex features and patterns from large amounts of data, allowing computer vision systems to better distinguish between true intrusions and normal changes in the environment. However, deep learning algorithms are slow, take up a lot of memory, and require the support of a high-performance computer. There are many cameras in cities, but fewer cameras are used for intrusion detection in urban rail transit, and it is not practical to use a large number of high-performance computers. An efficient method for urban rail intrusion detection is needed in complex transportation scenarios.

## II. LITERATURE REVIEW

Intrusion detection is an active research topic in the field of urban rail transit security. At present, some universities and research institutions in the United States are also conducting in-depth research in this area, such as the University of California, Berkeley, Yale University and so on. They have obtained more accurate and practical algorithmic models by introducing deep learning techniques and combining a large number of data sets for training, and target detection algorithms based on convolutional neural networks have come into being and have been gradually applied to a variety of urban environments, which are now broadly categorized into single-stage detection algorithms and two-stage detection algorithms. Common single-stage detection algorithms include RetinaNet [1], SSD [2] and YOLO [3] series. The algorithm regards localization and classification as a regression problem, realizes end-to-end detection, and has a faster detection speed, but its anchor mechanism based method [4], generates a large number of candidate frames in detection, and the number of candidate frames in which the target is detected is small. Common two-stage detection algorithms such as R-CNN [5], Fast R-CNN [6] have been widely used in this field. They aim to improve the detection performance by reducing the redundancy in the candidate frames generated by the anchor mechanism. Two-stage detection algorithms first screen out all positive samples and subsequently generate regions of interest (ROIs), and then in the second stage of these two-stage detection algorithms, the bounding frames generated in the previous stage are further

refined by performing region classification and positional adjustments on the regions of interest. The whole process requires iterative detection, classification and position refinement, and although two-stage detection algorithms tend to be slower compared to single-stage detection algorithms, they usually achieve higher detection accuracy. Zuoming [7] proposed a CNN model for road extraction that optimizes the extracted data to obtain comprehensive road features. Jiguang Dai [8] introduced a multi-scale CN N based road extraction method for remote sensing images. They used sub-image training model and combined it with residual linking to solve the resolution reduction and gradient vanishing problems in the extraction process. X. Zhang [9] developed a FCN network which utilizes a spatially consistent integration algorithm to determine the weights of the loss function used to extract the road regions. Xiangwen Kong [10] introduced a SM-Unet semantic segmentation network with a strip pooling module to enhance the road extraction performance. Hao Qi [11] proposed an MBv 2-DPPM model considering segmentation accuracy and speed. He and Ren [12] proposed a train obstacle detection method based on improved R-CNN. A new parallel upsampling structure and a context extraction module were added to the architecture to improve the accuracy of R-CNN to 90.6%. In order to enhance the ability to recognize small target objects, He [13] applied an enhanced Mask RCNN railroad obstacle detection method and proposed a new feature extraction network that incorporates a number of multiscale improvement techniques. The two-level target detection technique is the basis of the above method. Multiscale feature fusion by extending the sensitive domain and fusing shallow and deep features can improve the target recognition ability to some extent. However, real-time detection is not possible due to the area suggestion network. Zhang [14] proposed a high speed rail intruder detection algorithm based on YOLOv3 network. By improving this algorithm with FPN structure and extracting features with switchable hollow convolution, the false alarm rate of target detection is reduced, but the frames per second (FPS) is low, which does not satisfy the requirement of real-time detection. Literature [15] proposed a lightweight adaptive multi-scale feature fusion target detection network based on YOLOv3 to improve the performance of small target detection in complex environments, but there is a leakage detection phenomenon for occluded objects Literature [16] proposed a traffic sign recognition algorithm based on YOLOv5, comparing the detection results of the current common algorithms for the traffic signs. Overall, YOLOv5s algorithm performs better and faster in urban track detection. However, there are still some challenges, such as dealing with complex traffic situations in the city and enhancing the detection speed of the model, so the yolov5s algorithm model is chosen as the experimental base model and improved.

This paper proposes an improved method for foreign object intrusion detection in urban rail transit using YOLOv5s network as the base network. The method is designed to meet the need for real-time and accurate detection in urban rail transit environment. The YOLOv5s network model consists of several parts including Input, Backbone, Feature Fusion, Neck and Head. To enhance the input data, a mosaic [17] data enhancement method is used at the Input. This technique randomly selects four images from the image library and then rearranges them to generate a new image. By using the mosaic technique, the YOLOv5s network can benefit from the increased variation and diversity in the training data, which helps to improve its performance in object detection tasks; the proposed method also includes an improvement to the channel attention mechanism used in the CBAM [18] module. This improvement yields a new attention mechanism module called SCBAM. The SCBAM module aims to enhance the information interaction between different channels, further improving the network's ability to capture relevant features and improving the overall performance of tasks such as object detection. The proposed SCBAM convolutional attention module replaces the original C3 module, which effectively improves the model's ability to specifically characterize and recognize small target objects. Three detection layers are used in the YOLOv5s detection head [19], which are responsible for predicting large, medium and small targets, respectively. In addition, the network structure, output characterization method and boundary regression loss function of the YOLOv5s architecture were modified to improve its performance. These modifications aim to reduce the false detection and leakage rates in urban rail transit intrusion detection. By enhancing the feature extraction capability of the algorithm, the improved YOLOv5s network architecture provides more accurate and reliable detection results for long-range and small intrusion objects.

## III. CS-YOLO ALGORITHM

### A. YOLOv5s Structural Model

YOLOv5s target detection framework has a great performance advantage in large and medium-sized target detection, with real-time monitoring speed up to 140 fps and relatively high recognition accuracy. YOLOv5s consists of 3 main parts: backbone network, bottleneck layer and detection head, and the model structure is shown in Fig. 1.

*1) Backbone network:* The YOLOv5s algorithm combines the C3 and SPPF modules to enhance its performance. The C3 module focuses on reducing computation and increasing inference speed, optimizing the model's efficiency. On the other hand, the SPPF [20] module conducts multi-scale feature extraction on a single feature map, contributing to improved model accuracy. By leveraging these modules, the enhanced YOLOv5s algorithm achieves a balance between computational efficiency and accurate object detection. This makes it well-suited for various applications requiring real-time, precise recognition, and tracking of objects, thereby enhancing detection accuracy and speed for small and medium-sized objects on the track. C3 contains three standard convolutional layers and several bottleneck modules, the number of which is determined by the configuration file. The number of bottleneck modules is determined by the product of the n and depth_multiple parameters of YAML.C3 is the main module for learning the residual features and is divided into two branches, one using the specified bottleneck layer and three standard convolutional layers, and the other passing through only one of the basic convolutional modules, and finally the two branches are merged. After the Concat

operation, the activation function in the standard convolution module is SiLU. This is a function that applies the Sigmoid linear unit by *elements, which is* characterized as take-anywhere, continuous, and smooth, not a monotonic function, and is suitable for representing nonlinear features. The principle of the SPPF module is basically the same as that of spatial pyramid pooling [21] but uses a different design of pooling kernels. SPP uses 4 pooling kernels by default in YOLOv5s: 5×5, 9×9, 13×13, and 1×1. SPPF uses two pooling kernels by default in YOLOv5s: 5×5 and 1×1.Spatial pyramid pooling allows fusing feature maps at different scales and in the SPP layer Apply different pooling operations to multiple scales to capture information from different receptive fields. By pooling features at different scales, the SPP layer increases the perceptual field of the network, allowing it to efficiently process inputs of different sizes. This approach allows the network to process inputs of various sizes without the need to resize or crop the image to a specific size beforehand. The SPP layer enhances the flexibility and adaptability of the network in processing inputs of different resolutions.



Fig. 1. YOLOv5s algorithm structure diagram.

*2) Neck network:* The YOLOv5s algorithm contains two key components, the Feature Pyramid Network (FPN) [22] and the Path Aggregation Network (PAN) [23]. The FPN utilizes a feature pyramid structure to integrate high-level semantic information with low-level features. This allows semantic knowledge to be passed top-down within the network. By combining the high-resolution details of low-level features with the contextual understanding of high-level

features, FPN enhances the overall semantic representation of the network. On the other hand, PAN facilitates bottom-up transfer of localization information so that low-level information is propagated to higher levels in the network. PAN achieves this by fusing feature information from feature maps of different sizes. By integrating FPN and PAN, the YOLOv5s algorithm benefits from the transfer of semantic information and the effective utilization of multi-scale features. By integrating FPN and PAN, the YOLOv5s algorithm benefits from both the delivery of semantic information and the effective utilization of multi-scale features.

*3) Head networks:* The head networks as the detection part of the model and is mainly used to predict objects of different sizes from the extracted multi-scale feature maps. The output anchor frame mechanism extracts a priori frame scales by clustering and constrains the location of the prediction boundaries. The first is an 8-fold downsampled output with respect to the input image, which has a small perceptual range, preserves low-level high-resolution features, and facilitates the detection of small objects. The second is a 16-fold downsampled output with respect to the input image, which has a medium perceptual range and is suitable for detecting medium objects. The third is an output downsampled 32 times with respect to the input image, which has a large perceptual range and is suitable for detecting large objects.

### B. Improved CBAM-based Attention Mechanism Module

The CBAM module is an attention mechanism that enhances spatial [24] and channel [25] attention while minimizing parameters. Its performance in classification detection on public datasets has shown improvement. However, when applied to the urban railroad track intrusion detection dataset in this paper, challenges arise. These challenges include a greater variety of vehicle classifications, a scarcity of samples, and the presence of vehicles with similar features. Consequently, the classification detection performance falls short of achieving the anticipated results. Fig. 2 illustrates the structure of the CBAM module.

The original channel attention mechanism employed global pooling operations, which included average pooling and maximum pooling on the input feature map along the width and height dimensions. These pooling operations aggregated information from all spatial locations within each channel. Following pooling, the generated features underwent processing through a multilayer perceptron (MLP). The MLP performed element-wise operations [26] to learn the importance weights for each channel. Finally, the output of the MLP passed through the Sigmoid [27] activation function to generate the final channel attention feature map. The Sigmoid function maps the values to a range between 0 and 1, representing the importance or activation level of each channel. In summary, the original channel attention mechanism comprises a pooling operation, an MLP for weighting, and a Sigmoid operation that produces the final channel attention feature map.

Fig. 2. CBAM structure diagram.



Fig. 3. SCAM structure diagram.

When information in the channel is globally and maximally pooled, this operation tends to neglect information interaction within the channel. In this paper, this paper optimizes and enhances the channel attention component of CBAM by discarding the maximum pooling operation. Instead, this paper represent the compressed features by aggregating the spatial information of the feature map through the average pooling operation. Drawing inspiration from the concept of ECA [28], this paper replaces the MLP network with a one-dimensional convolution operation in the improved version of the channel attention mechanism. Following one-dimensional convolutional processing, the ability to interact information between different channels is strengthened, and the output is subsequently simplified by a sigmoid function. In summary, as illustrated in Fig. 3, the improved channel attention module (SCAM) achieves local cross-channel interaction information of size K. This enhancement boosts the model's feature extraction capability, particularly for small and medium-sized objects.

*C. KM-Decoupled Head*

When recognizing targets on urban tracks, the task becomes more challenging due to the presence of multiple types of targets. To address this challenge and enhance the accuracy of localization and classification, this paper introduces a decoupled detection head known as KM-Decoupled Head. The primary idea behind this approach is to separate the feature channels for localization and classification tasks, specifically for bounding box coordinate regression and object classification. This decoupling enables a more precise estimation of bounding box coordinates, ultimately improving localization accuracy. Additionally, it ensures that the features used for classification are less influenced by changes in the localization task, leading to improved object classification. In essence, the goal of KM-Decoupled Head is to enhance target prediction in scenarios involving multiple types of targets and occluded targets. By decoupling the feature channels used for localization and classification, it enhances the accuracy of both tasks, resulting in more efficient target identification on urban tracks.

As illustrated in Fig. 4, the KM-Decoupled Head, a decoupled detection head, follows a specific process. First, a 1×1 convolution [29] is applied to the input feature map to reduce the number of channels and parameters generated. Subsequently, the output feature map is split into two branches to address the classification and localization tasks separately. For the classification branch, features are extracted using a 3×3 deep convolution [30]. The number of channel bits in the feature map is then adjusted to match the predicted number of target categories through a 1×1 convolution. On the other hand, the localization branch also employs a 3×3 deep convolution for feature extraction. After feature extraction, the resulting feature map is divided into two parts. One part predicts the center coordinates of the bounding box, along with the height and width of the bounding box. In summary, the KM-Decoupled Head employs a 1×1 convolution to reduce the number of channels and subsequently splits the feature map into separate branches for classification and localization tasks. The classification branch adjusts the feature map channels based on the predicted target categories, while the localization branch predicts the bounding box coordinates and confidence scores for target identification.

The decoupled structure of the KM-Decoupled Head offers several advantages over the coupled detector head, which integrates multiple types of information into a single feature map. Firstly, the decoupled design effectively avoids potential conflicts between different task requirements, thereby enhancing localization and classification capabilities. By separating the feature channels for each task, the model can focus on learning distinct representations for precise localization and accurate classification. Secondly, the decoupling head preserves information in each channel through depth and breadth operations, ensuring that valuable information is not weakened or diluted during processing steps. Additionally, this approach helps reduce computational complexity, thereby accelerating network convergence. Finally, due to the depth and breadth operations, the decoupling head achieves faster inference. In conclusion, the decoupled structure of the KM-Decoupled Head improves localization and classification by mitigating conflicts, effectively preserving channel information, reducing computational complexity, and achieving faster inference speed through depth and breadth operations.

**Couple Head**



Fig. 4.   Coupling head with KM-decoupled head.

### D.  Loss Function SIoU

The loss function in the YOLOv5s model is CIoU[31] :

$$L_{CIoU} = 1 - IoU + \frac{(b,\ b^{gt})p^2}{c^2} + \alpha v \qquad (1)$$

$$v = \frac{4}{\pi^2}\left(\arctan\frac{w^{gt}}{h^{gt}} - \arctan\frac{w}{h}\right)^2 \qquad (2)$$

$$\alpha = \frac{v}{(1-IoU)+v} \qquad (3)$$

where IoU is the ratio of the intersection and union of the predicted and actual frames, b represents the center of the predicted point, $b^{gt}$ represents the center of the actual frame, p denotes the Euclidean distance, c denotes the length of the diagonal of the bounding box formed by the predicted frame and the actual frame, $\alpha$ represents the weighting coefficient, and v represents the difference in the aspect ratio of the predicted frame and the actual frame.

CIoU does not consider the direction of mismatch between the actual frames and the predicted frames, leading to slower convergence and lower efficiency. Therefore, we introduce a more balanced loss function, SIoU [32], which increases the vector angle between regressions. We also redefine the cost function (penalty indicator), effectively reducing the degrees of freedom of regression, speeding up the convergence of the network, and further improving the accuracy of regression. SIoU consists of four cost functions.

- Angle_Loss Angle_Cost minimizes the number of variables associated with distance. The formula is:

$$\lambda = 1 - 2 * \sin^2\left(\arcsin\left(\frac{c_h}{\sigma}\right) - \frac{\pi}{4}\right) \qquad (4)$$

where $\sigma$ is the distance between the center point of the real frame and the predicted frame, and ch is the height difference between the center point of the real frame and the predicted frame.

- Distance_LossDistance_Cost explores the distance of different bounding boxes at different centers as much as possible. The formula is:

$$\Delta = \sum_{t=x,\ y}(1 - e^{-\gamma p_t}) \qquad (5)$$

$$p_x = \left(\frac{b_{c_x}^{gt}-b_{c_x}}{c_w}\right)^2, \quad p_y = \left(\frac{b_{c_y}^{gt}-b_{c_y}}{c_h}\right)^2 \qquad (6)$$

"cw" and "ch" refer to the width and height of the minimum bounding rectangle of the actual box and the predicted box, respectively. "$b_{c_x}^{gt}$" and "$b_{c_y}^{gt}$" represent the coordinates of the center of the actual box, while "$b_{c_x}$" and "$b_{c_y}$" represent the center coordinates of the predicted frame.

- The shape loss Shape_Cost represents the deviation of the center of the predicted frame from the center of the real frame in an effort to obtain the optimal predicted frame. The formula is:

$$\Omega = \sum_{t=w,\ h}(1 - e^{-w_t})^\theta \qquad (7)$$

$$\omega_w = \frac{|w-w^{gt}|}{max(w,\ w^{gt})} \qquad (8)$$

$$\omega_h = \frac{|h-h^{gt}|}{max(h,\ h^{gt})} \qquad (9)$$

"w" and "h" represent the width and height of the predicted box, while "wgt" and "hgt" represent the width and height of the actual box,$\theta$ is the degree of concern for shape loss.

- Iou_Cost is the ratio of the intersection and the concatenation between the predicted and real boxes. The formula is:

$$IoU = \frac{|B \cap B^{GT}|}{|B \cup B^{GT}|} \qquad (10)$$

- the regression loss function SIoU is:

$$Loss_{SIoU} = 1 - IoU + \frac{\Delta + \Omega}{2} \qquad (11)$$



Fig. 5. Improved YOLOv5s backbone network.

### E. Backbone Network Improvement

The enhanced SCBAM attention mechanism module replaces the original C3 module in the YOLOv5s backbone. It strengthens inter-channel information fusion using the channel attention mechanism and improves inter-channel information

fusion through the spatial attention mechanism. The spatial attention mechanism is focused on detecting the target's location. By combining these two mechanisms, the output information becomes more concentrated on key features, enabling the model to prioritize important features during target detection. This enhances the feature extraction capability and ultimately improves the model's accuracy in detecting objects of various sizes. Additionally, this paper replaces the Focus slicing operation with a 6x6 convolution and substitutes the SPP module with the SPPF module to address the speed issue caused by parallel operations in the original model. The enhanced backbone network diagram is illustrated in Fig. 5

### IV. EXPERIMENTAL ANALYSIS

#### A. Experimental Environment and Parameter Settings

The specific experimental environment is shown in Table I.

TABLE I. EXPERIMENTAL ENVIRONMENT

| Configuration name | Version/parameters |
|---|---|
| Operating system | Windows 11 |
| Video storage | 16GB |
| GPU | RTX4050 |
| Memory | 64GB |
| Python | 3.7 |
| Deep learning framework | pytorch1.11.0 |

During training, the network is trained using the SGD optimizer with an initial learning rate of 0.2. The learning rate is adjusted using the cosine annealing strategy. The batch size is set to 16 and the training process is performed for a total of 300 epochs.

#### B. Experimental Data Set

To assess the superiority of the improved target detection algorithm proposed in this paper, experiments were conducted on urban rail foreign object intrusion using the custom dataset, WZ-dataset.

As there is no publicly available dataset for foreign objects intruding on tracks in cities. This article develops a new dataset WZ-dataset. This dataset comprises 10 common types of foreign objects that may intrude on urban rail tracks, including pedestrians (person), bicycles (bike), electric bicycles (ebike), three-wheeled vehicles (tricycle), bottles, bags, cars, buses, trucks, and books. The WZ-dataset contains a total of 2,000 images, distributed across training, validation, and testing sets in an 8:1:1 ratio. Specifically, there are 1,600 images in the training set, 200 images in the validation set, and 200 images in the test set. The image size is consistently maintained around 600×800 pixels. To ensure a balanced sample size, each type of target is evenly distributed. The CS-YOLO algorithm model proposed in this paper is trained using the training and validation sets of the WZ-dataset (see Fig. 6). The model's performance is then evaluated by measuring the final average accuracy and detection speed on the test set.

Fig. 6. Example of WZ dataset.

## C. Performance Index

The performance metrics, including Precision (P), Recall (R), Mean Accuracy (mAP), Parameters, Inference Time per Image, and Frames per Second (FPS), are utilized in the experiments to evaluate the proposed algorithm's performance. The calculation formula for these performance metrics is as follows:

$$\text{Precision} = \frac{TP}{TP+FP} \qquad (12)$$

$$\text{Recall} = \frac{TP}{TP+FN} \qquad (13)$$

$$mAP = \frac{\sum_{i=1}^{N} AP_i}{N} \quad AP = \int_0^1 P(R)dR \qquad (14)$$

$$FPS = \frac{TotalTime}{NumFigure} \qquad (15)$$

When evaluating object detection models, TP refers to true positives, representing instances where the model predicts correctly. FP stands for false positives, indicating instances where the model predicts incorrectly. Average Precision (AP) measures the area under the precision-recall curve, offering a comprehensive evaluation of the model's performance at various thresholds concerning both precision and recall. AP calculates the average precision for each category, and Mean Average Precision (mAP) is the average of the APs across all categories. mAP provides an overall assessment of the model's performance across multiple object categories. mAP0.5 represents the average precision calculated at an IoU threshold of 0.5, considering a bounding box prediction correct if the Intersection over Union (IoU) is greater than or equal to 0.5. For mAP0.5:0.95, average accuracy is calculated across IoU thresholds ranging from 0.5 to 0.95 in steps of 0.05. Frames per Second (FPS) represents the number of images processed and detected by the network model per second. It reflects the speed and efficiency with which the network model performs object detection. A higher FPS value indicates that the model can process images more quickly, enabling real-time or near real-time applications.

## D. Experimental Analysis of SCBAM Attention Module

The SCBAM (Spatial and Channel Bottleneck Attention Module) is an enhanced version of the CBAM (Convolutional Block Attention Module). This paper aims to assess the effectiveness of the SCBAM attention module through five sets of comparative experiments on the WZ dataset. The experimental results presented in Table II demonstrate the efficacy of the SCBAM attention module. On the WZ dataset, the average accuracy mAP0.5 improved by 3.4%, and mAP0.5:.95 improved by 6.3%. Importantly, these enhancements were achieved while maintaining good detection speed, ensuring real-time detection performance. These findings strongly support the effectiveness of the SCBAM attention module proposed in this paper.

TABLE II.    COMPARISON TESTS OF DIFFERENT ATTENTION MECHANISMS ON THE WZ DATASET

| Method | mAP@0.5 | mAP@0.5:.95 | FPS |
|---|---|---|---|
| Baseline | 0.844 | 0.532 | 85 |
| +SENet | 0.856(+1.2)% | 0.556 (+2.4)% | 77 |
| +CA | 0.865 (+2.1)% | 0.505 (-2.7)% | 90 |
| +ECA | 0.867 (+2.3)% | 0.574 (+4.2)% | 73 |
| +CBAM | 0.872 (+2.8)% | 0.587 (+5.5)% | 83 |
| +SCBAM | 0.878 (+3.4)% | 0.595 (+6.3)% | 91 |

*E. Ablation Experiments*

CS-YOLO is an enhanced version of the YOLOv5s network model, featuring improvements in various aspects, including the attention mechanism, output characterization method, backbone network, and bounding box regression loss function. In this paper, a series of ablation experiments are conducted to assess the impact of the proposed modules on the algorithm's performance. These experiments analyze the performance optimization achieved by individual modules, as well as the combined effect when different modules are introduced in different orders. The objective is to determine the effectiveness and contribution of each module in enhancing the overall performance of the algorithm. The results of the ablation experiments are presented in Table III.

This paper introduces three improved methods incorporated into the YOLOv5s network model (indicated by a checkmark in Table III). These methods enhance the detection accuracy on the WZ dataset to varying degrees. The CS-YOLO algorithm proposed in this paper achieves a detection speed of 78 FPS, outperforming the original YOLOv5s algorithm. On the WZ dataset, this algorithm demonstrates a 4.9% improvement in mAP0.5 and an 8.9% improvement in mAP0.5:95. Experiments conducted on the WZ dataset validate the effectiveness of the proposed algorithm in urban rail transit detection. These experiments showcase the algorithm's capability to address the challenge of foreign object intrusion detection in complex urban rail transit environments.

*F. Comparison Experiments*

Table IV in this paper presents experimental results comparing the CS-YOLO algorithm with several other existing algorithms on the WZ dataset. The compared algorithms include the original YOLOv5s, SSD, Faster R-CNN, YOLOv3, YOLOv4-tiny, YOLOv4, YOLOv5m, YOLOX-tiny, YOLOX-S, YOLOv6-tiny, and YOLOv7-tiny. The experimental results comprehensively analyze the performance of the CS-YOLO algorithm in comparison with these popular algorithms. For detailed information about the experimental results and algorithm performance comparison on the RS dataset, please refer to Table IV in this paper.

The experimental results in Table IV demonstrate that the CS-YOLO algorithm proposed in this paper achieves higher detection accuracy on the WZ dataset compared to other mainstream algorithmic models. Notably, the parameters of the algorithm proposed in this paper are reduced by 130%, while the detection accuracy is improved by 2.7% compared to the YOLOv5m network with a similar structure. Although YOLOV4-tiny and YOLOv7-tiny exhibit higher detection speeds (111 FPS and 119 FPS, respectively), their detection accuracies are relatively low at 75.5% and 83.3%, making them unsuitable for complex urban rail transit environments. In conclusion, the CS-YOLO algorithm proposed in this paper demonstrates the highest detection accuracy while maintaining good real-time performance, offering a significant advantage over other algorithms under consideration.

TABLE III.       CS-YOLO ABLATION EXPERIMENTS ON THE WZ DATASET

| Group | SCBAM | KM-DHead | SIoU | mAP@0.5 | mAP@0.5:.95 | Parameters | FPS |
|---|---|---|---|---|---|---|---|
| 1 | | | | 0.844 | 0.532 | 7.1M | 85 |
| 2 | √ | | | 0.878 (+3.4)% | 0.595 (+6.3)% | 7.1M | 91 |
| 3 | | √ | | 0.878 (+3.4)% | 0.553 (+2.1)% | 7.4M | 74 |
| 4 | | | √ | 0.848 (+0.4)% | 0.558 (+2.6)% | 7.1M | 94 |
| 5 | √ | √ | | 0.884 (+4.0)% | 0.590 (+5.8)% | 7.4M | 85 |
| 6 | | √ | √ | 0.878 (+3.4)% | 0.592 (+6.0)% | 7.4M | 80 |
| 7 | √ | | √ | 0.881 (+3.7)% | 0.589 (+5.7)% | 7.1M | 78 |
| 8 | √ | √ | √ | 0.893 (+4.9)% | 0.621 (+8.9)% | 7.4M | 78 |

TABLE IV.       COMPARATIVE EXPERIMENTS OF COMMON ALGORITHMS ON THE WZ DATASET

| Method | mAP@0.5(%) | mAP@0.5:.95(%) | Parameters(M) | Inference(ms) | FPS |
|---|---|---|---|---|---|
| YOLOv5s | 0.844 | 0.532 | 7.1 | 11.3 | 85 |
| SSD | 0.865 | 0.639 | 100.2 | 23.1 | 55 |
| Faster R-CNN | 0.711 | 0.453 | 136.2 | 46.6 | 23 |
| YOLOv3-spp | 0.833 | 0.512 | 9.56 | 12.3 | 81 |
| YOLOv4-tiny | 0.755 | 0.502 | 5.9 | 10.9 | 111 |
| YOLOv4 | 0.817 | 0.528 | 244.8 | 55.2 | 35 |
| YOLOv5m | 0.866 | 0.636 | 20.9 | 19.9 | 70 |
| YOLOX-tiny | 0.822 | 0.565 | 5.1 | 16.8 | 43 |
| YOLOX-S | 0.848 | 0.591 | 9.0 | 23.5 | 25 |
| YOLOv6-tiny | 0.866 | 0.611 | 15.0 | 22.5 | 81 |
| YOLOv7-tiny | 0.833 | 0.567 | 6.1 | 8.3 | 119 |
| CS-YOLO | 0.893 | 0.621 | 7.4 | 14.8 | 78 |

In this paper, the detection performance before and after the improvement is visually compared, as illustrated in Fig. 7. The comparison results demonstrate that the CS-YOLO algorithm exhibits excellent detection performance in both sets of images. The results further highlight that the CS-YOLO algorithm effectively addresses misdetection and omission issues when detecting fuzzy and small targets in complex urban rail backgrounds compared to the initial YOLOv5s algorithm. With an FPS of 78, the CS-YOLO algorithm achieves higher detection accuracy by more accurately identifying common foreign object features while maintaining real-time detection speed. This capability meets the need for real-time and accurate detection in complex urban rail scenes.

### G. Analysis of Experimental Results

The accuracy recall curve of the initial YOLOv5s and CS-YOLO are shown in Fig. 8. The horizontal axis represents the recall rate, and the vertical axis represents precision. The inset box displays the detection precision for various common intruders, with the bolded blue line representing the average precision across all detection categories. Upon comparing the detection results of the two algorithms, it is evident that the CS-YOLO algorithm exhibits lower detection precision for tricycles and trucks compared to the initial YOLOv5s. However, it demonstrates higher detection precision for all other intrusions. The average accuracy is improved by 4.9% compared to the initial YOLOv5s algorithm, indicating the significance of the proposed improvements.



| Original detection chart | YOLOv5s detection effect | CS-YOLO detection effect |

Fig. 7.   YOLOv5s and CS-YOLO detection effect comparison.



| Graph of CS-YOLO detection results | YOLOv5s detection result chart |

Fig. 8.   Graph of CS-YOLO and YOLOv5s detection results.

## V. CONCLUSION

In this paper, we propose the CS-YOLO algorithm model for detecting foreign object intrusions in complex urban railways. The algorithm addresses the challenges of low accuracy and poor timeliness present in existing methods. Key enhancements include the introduction of the SCBAM attention mechanism module, replacing the original C3 module in YOLOv5s. Additionally, the algorithm features the KM-Decoupled Head decoupling head, utilizes SIoU as the loss function, and modifies the backbone network structure. In comparison to other mainstream target detection algorithms, This article improves the algorithm achieves superior detection accuracy, particularly in resolving issues of false positives and missed detections when dealing with concealed and small targets. Despite the enhanced accuracy, the algorithm maintains real-time detection speed, making it well-suited for detecting foreign objects on complex urban tracks. Future research should focus on optimizing the network for easier deployment on embedded GPU platforms, ensuring its applicability in real-world scenarios.

## REFERENCES

[1] Qiu Zhiruo,Rong Shiyang,Ye Likun. YOLF-ShipPnet: Improved RetinaNet with Pyramid Vision Transformer[J]. International Journal of Computational Intelligence Systems,2023,16(1).

[2] Chen Renfei,Wu Jian,Peng Yong,Li Zhongwen,Shang Hua. Solving floating pollution with deep learning: a novel SSD for floating objects based on continual unsupervised domain adaptation[J]. Engineering Applications of Artificial Intelligence,2023,120.

[3] Li Jie,Li Sudong,Li Xiaoli,Miao Sheng,Dong Cheng,Gao Chuanping,Liu Xuejun,Hao Dapeng,Xu Wenjian,Huang Mingqian,Cui Jiufa. primary bone tumor detection and classification in full-field bone radiographs via YOLO deep learning model.[J]. European radiology,2023,33(6).

[4] Said Yahia,Atri Mohamed,Albahar Marwan Ali,Ben Atitallah Ahmed,Alsariera Yazan Ahmad. Indoor Signs Detection for Visually Impaired People. Navigation Assistance Based on a Lightweight Anchor-Free Object Detector[J]. International Journal of Environmental Research and Public Health,2023,20(6).

[5] Zhang Wenming,Zhu Qikai,Li Yaqian,Li Haibin. MAM Faster R-CNN: Improved Faster R-CNN based on Malformed Attention Module for object detection on X-ray security inspection[J]. Digital Signal Processing,2023,139.

[6] Sang-Soo Baek,JongCheol Pyo,Yakov Pachepsky,Yongeun Park,Mayzonee Ligaray,Chi-Yong Ahn,Young-Hyo Kim,Jong Ahn Chun,Kyung Hwa Cho. Identification and enumeration of cyanobacteria species using a deep neural network[J]. Ecological Indicators,2020,115(C).

[7] She, Z.; Shen, Y.; Song, J.; Xiang, Y. Using the classical CNN network method to construct the automatic extraction model of remote sensing image of Guiyang road elements. Bull. Surv. Mapp. 2023, 4, 177‑182.

[8] Dai, J.; Du, Y.; Jin, G.; Tao, D. A Road Extraction Method Based on Multiscale Convolutional Neural Network. Remote Sens. Inf. 2019, 35, 28‑37.

[9] Zhang, X.; Ma, W.; Li, C.; Wu, J.; Tang, X.; Jiao, L. Fully Convolutional Network-Based Ensemble Method for Road Extraction from Aerial Images. IEEE Geosci. Remote Sens. Lett. 2019, 17, 1777‑1781.

[10] Kong, X.; Wang, C.; Zhang, S.; Li, J.; Sui, Y. Application of Improved U-Net Network in Road Extraction from Remote Sensing Images. Remote Sens. Inf. 2022, 37, 97‑104.

[11] Qi, H.; Li, Y.; Qi, Y.; Liu, L.; Dong, Z.; Du, X. Research on Track and Obstacle Detection Based on New Lightweight Semantic Segmentation Network. J. Railw. Sci. 2019, 45, 58‑66.

[12] He, D.; Ren, R.; Li, K.; Zou, Z.; Ma, R.; Qin, Y.; Yang, W. Urban Rail Transit Obstacle Detection Based on Improved R-CNN. Measurement 2022, 196, 111277.

[13] He, D.; Qiu, Y.; Miao, J.; Zou, Z.; Li, K.; Ren, C.; Shen, G. Improved Mask R-CNN for Obstacle Detection of Rail Transit.Measurement 2022, 190, 110728.

[14] Zhang, J.; Wang, D.; Mo, G. High-speed rail foreign body intrusion detection algorithm based on improved YOLOv3. Comput.Technol. Dev. 2022, 32, 69‑74.

[15] Ye, T.; Zhao, Z.; Wang, S.; Zhou, F.; Gao, X. A Stable Lightweight and Adaptive Feature Enhanced Convolution Neural Network for Efficient Railway Transit Object Detection. IEEE Trans. Intell. Transport. Syst. 2022, 23, 17952‑17965.

[16] Yingning Gao , Weisheng Liu. Complex Labels Text Detection Algorithm Based on Improved YOLOv5[J]. IAENG International Journal of Computer Science,2023,50(2).

[17] Mech Mario,Ehrlich André,Herber Andreas,Lüpkes Christof,Wendisch Manfred,Becker Sebastian,Boose Yvonne,Chechin Dmitry,Crewell Susanne,Dupuy Régis,Gourbeyre Christophe,Hartmann Jörg,Jäkel Evelyn,Jourdan Olivier,Kliesch LeifLeonard,Klingebiel Marcus,Kulla Birte Solveig,Mioche Guillaume,Moser Manuel,Risse Nils,RuizDonoso Elena,Schäfer Michael,Stapf Johannes,Voigt Christiane. author Correction: MOSAiC-ACA and AFLUX - Arctic airborne campaigns characterizing the exit area of MOSAiC.[J]. Scientific data,2023,10(1).

[18] Liu Hui,Yang Guangqi,Deng Fengliang,Qian Yurong,Fan Yingying. MCBAM-GAN: The Gan Spatiotemporal Fusion Model Based on Multiscale and CBAM for Remote Sensing Images[J]. Sensing Images[J]. Remote Sensing,2023,15(6).

[19] Cao Lianyu,Zhang Xiaolu,Wang Zhaoshun,Ding Guangyu. Multi Angle Rotation Object Detection for Remote Sensing Image Based on Modified Feature Pyramid Networks[J]. International Journal of Remote Sensing,2021,42(14).

[20] Buffington Matthew L., Garretson Alexis, Kula Robert R., Gates Michael W., Carpenter Ryan, Smith David R., Kula Abigail A.R.. Pan trap color preference across Hymenoptera in a forest clearing[J]. Entomologia Experimentalis et Applicata,2020,169(3).

[21] Kase Kina,Nakanishi Yosuke,Murono Shigeyuki,Yoshizaki Tomokazu. Comparison of Plain and Contrast-enhanced Computed Tomography for the Detection Head-and-neck Abscess[J]. Practica oto-rhino-laryngologica. Suppl.,2018,152(0).

[22] Nishiyama T.,Kumagai A.,Kamiya K.,Takahashi K.. SILU: Strategy Involving Large-scale Unlabeled Logs for Improving Malware Detector[J]. Proceedings - IEEE Symposium on Computers and Communications,2020,2020-July.

[23] Nagaoka Hiroshi,Ohtake Makiko,Karouji Yuzuru,Kayama Masahiro,Ishihara Yoshiaki,Yamamoto Satoru,Sakai Risa. Sample studies and SELENE (Kaguya) observations of purest anorthosite (PAN) in the primordial lunar crust for future sample return mission[J]. Icarus,2023,392.

[24] Yu Junwei,Shen Yi,Liu Nan,Pan Quan. Frequency-Enhanced Channel-Spatial Attention Module for Grain Pests Classification[J]. Agriculture,2022,12(12).

[25] Lee Haeyun,Cho Sunghyun. Image Restoration Network with Adaptive Channel Attention Modules for Combined Distortions[J]. Journal of the Korea Computer Graphics Society,2019,25(3).

[26] Kim Dohyun,Park Daeyoung. Element-Wise Adaptive Thresholds for Learned Iterative Shrinkage Thresholding Algorithms[J]. IEEE Access,2020,8.

[27] Coronavirus - COVID-19; Hanchuan People's Hospital Reports Findings in COVID-19 (The Challenges Of Urgent Radical Sigmoid Colorectal

Cancer Resection In A COVID-19 Patient; A Case Report)[J]. Medical Letter on the CDC & FDA,2020.

[28] Wang Xinkai,Jia Xu,Zhang Miyuan,Lu Houda. Object Detection in 3D Point Cloud Based on ECA Mechanism[J]. Journal of Circuits, Systems and Computers,2023,32(05).

[29] Yang K. Research on the Application of Time Convolution Series in Futures Price Forecasting[C]//Wuhan Zhicheng Times Cultural Development Co. ...Proceedings of 3rd International Conference on the Frontiers of Innovative Economics and Management (FIEM 2022).BCP Business & Management. BCP Business & Management, 2022:151-155.DOI:10.26914/c.cnkihy.2022.069652.

[30] Donghan X,Zhi W,Chunlin C, et al. Depthwise Convolution for Multi-Agent Communication With Enhanced Mean-Field Approximation.[J]. IEEE transactions on neural networks and learning systems,2022,PP.

[31] Yifan J,Dexin G,Shiyu Z, et al. A real-time fire detection method from video for electric vehicle-charging stations based on improved YOLOX-tiny [J]. . Journal of Real-Time Image Processing,2023,20(3).

[32] Yonghong L,Cheng Z,Zhiqiang Z, et al. Research on detection method of Tubercle Bacilli based on the improved YOLOv5.[J]. Physics in medicine and biology,2023,68(10).

# Students' Perception of ChatGPT Usage in Education

Irena Valova, Tsvetelina Mladenova, Gabriel Kanev
Computer Systems and Technologies, University of Ruse, Ruse, Bulgaria

*Abstract*—This research article delves into the impact of ChatGPT on education, focusing on the perceptions and usage patterns among high school and university students. The article begins by introducing ChatGPT, emphasizing its rapid user adoption and widespread interest. It explores the application of ChatGPT in various fields, including healthcare, agriculture, and education. A comprehensive survey involving 102 students, both high school and university, is detailed, covering aspects like familiarity with ChatGPT, reasons for usage, self-assessment of its effectiveness, and attitudes toward informing teachers about its use. The findings reveal varied perspectives on the benefits and challenges of incorporating ChatGPT in the learning process. The article concludes by emphasizing the need for careful consideration and integration of AI technologies in education, highlighting the risks of uncritical reliance on such tools and advocating for a balanced approach to foster students' critical thinking and intellectual growth.

*Keywords—Artificial intelligence in education; assessment; ChatGPT; Generative Pretrained Transformer 3, GPT-3; higher education; learning, teaching; Natural Language Processing (NLP)*

## I. INTRODUCTION

ChatGPT is a language model (chatbot) created by OpenAI that allows humans to interact with a computer naturally. A chatbot is an application used to conduct a conversation through the exchange of text messages or text-to-speech between a human and a computer/machine. These are computer programs that can hold a conversation with a user in natural language, understand their intent, and respond based on predefined rules and data. Designed to convincingly simulate the way a human would behave as a conversation partner, chatbot applications typically require constant tuning and testing. While working they are self-educating and improving.

For many researchers and for high education itself, it is important to see how high-school students (in their final years of high school) and university students perceive the idea of using such chatbots in their studies. This article is an examination of the adoption of AI by the students – how they are using it, how frequently, what type of questions they ask it, to what degree they understand the answers, and how they implement them in their class assignments.

This article is the result of an analysis of a questionnaire given among 102 Bulgarian students. It presents the questions, their answers, and some thoughts about the results. While the survey was anonymous, the respondents are students, the authors, are teaching and therefore we have first-hand observations about their problem-solving skills and their thought patterns.

It is obvious that this type of AI is here to stay, and it is up to the universities how they will be able to adopt and use it. Conducting such surveys will help them to understand it better and apply it efficiently.

## II. LITERATURE REVIEW

OpenAI is an artificial intelligence (AI) research and Implementation Company ensuring that general-purpose artificial intelligence benefits all of humanity. The company is dedicated to putting this alignment of interests first even before profit.

The definition of AI characterizes it as a branch of computer science that deals with the automation of intelligent behavior. The degree of intelligence is difficult to define, and therefore artificial intelligence cannot be precisely defined either. The term is used to describe systems that aim to use machines to emulate and simulate human intelligence and related behavior. This can be achieved through simple algorithms and predefined models, but it can also become much more complex.

ChatGPT (Chat Generative Pre-trained Transformer) was publicly presented in the summer of 2020 and launched in November 2022. It is an object of curiosity, controversy, and scientific interest among a wide range of Internet users from all ages and stages of life. Unlike search engines (such as Google, Bing, or Baidu), ChatGPT does not crawl the web for information about current events and information, and its knowledge is limited to things it learned before January 2022. It is the subject of many comments and discussions, from the fact that some analysts see it as a threat to some professions, to the fact that others believe that this technology is extremely successful and useful. Although this is not the first application based on artificial intelligence, it can be said that it is the most tested and has generated the most interest among users. The first million users were reached in just five days, which for other platforms took months and years (for example, Facebook reached a million users in 10 months in 2004). In just three months, ChatGPT users reached one billion (see Fig. 1), [https://www.tooltester.com/en/blog/ChatGPT-statistics/]. The first reactions are obviously to test and see how this brand-new technology works, and if it works. Almost immediately after testing, reasonable questions arise as to how useful and how dangerous such technology is. Many studies have analyzed the impact on different professions and business fields [1], and the impact on different fields of study [2] and industry [3].

In study [4], the authors conclude that the presence of various AI agents such as ChatGPT will change the context of higher education, but this will not be disruptive. It is very important to realize and assess this transformation promptly and to model it appropriately.

Fig. 1. ChatGPT visitors since release.

The impact of ChatGPT is examined in various fields, such as healthcare, medicine, and dentistry [5, 6]. Its impact and expectations in the field of agriculture and livestock breeding are also examined.

Today's agriculture uses various smart technologies and collects large volumes of data that can be used for crop forecasting, soil analysis, identification of crop diseases and pests, precision farming and irrigation planning, animal behavior analysis, and assessment of their condition. This can be helpful for businesses to make informed decisions and increase their profits. In study [7], the author investigates the potential positives and negatives of the application of chat GPT in agriculture. He provides examples of questions where ChatGPT can be useful in agriculture by analyzing and editing its answers. These are assessments of atmospheric conditions, soil, and air quality, diseases of different crops, and others. The author is noting the following points of ChatGPT usage that are valid for every other usage area:

- Strong dependence on data quality - if the data is inaccurate, biased, or incomplete - this will inevitably affect the responses from the agent;

- Lack of experience - ChatGPT is good at analyzing data but is not a specialist in any specific field and it is very important to have an experienced professional in the relevant field to interpret the model results and make the final decisions. The human factor cannot be avoided, and the specialists have the final say on the decisions.

Will it replace university professors or classroom teachers - this question is being asked more and more often. In research [8], the authors made a qualitative analysis based on a methodology for data collection, documentation, and drawing conclusions, which analysis shows that ChatGPT can only be a tool in training, and it is not possible to completely replace the trainer. It is more important to find an appropriate way to integrate technology into the learning process and, respectively, to develop the competencies of teachers in managing learning with such technologies.

Some researchers evaluate ChatGPT as an opportunity to increase the effectiveness of training and the motivation of students [9, 10, and 11] because the use of this new technology allows learning at an individual pace of the student and because students choose the direction of deepening their knowledge, they are much more motivated. ChatGPT provides personalized and interactive help, which engages more learners and develops self-learning skills.

It is natural to think in the direction of whether to allow or block access to ChatGPT in educational institutions [12] or whether to look for applications that detect if a given text is generated by ChatGPT. Before any measurements are taken it is interesting to see whether it is used and to what extent is used by the students.

The detailed analysis of the possibilities and limitations of ChatGPT made in [13, 14] shows that the use of this technology has great potential for application in the field of education, but it also comes with quite a few limitations and challenges. The use of ChatGPT, with all its positives and negatives, in training is in its very early stages and this assumes much more research in this area.

The results of the analysis of the use of ChatGPT in the different areas of higher education show a major problem in scientific writing [15]. Is it plagiarism to use text generated by ChatGPT and how to reference such text, what percentage of such text, relative to the total volume, is permissible in student and faculty scholarly publications are some of the topics examined by the author.

The study in [16] concluded that in the context of using ChatGPT in education, it should be noted that technology can only be a tool and cannot completely replace the role of the teacher. Therefore, it is necessary to integrate technology into learning appropriately and effectively and to develop the competence of teachers in managing learning with such technologies.

How can the use of such technologies be useful?

- Can be used to search for information and ask questions from different fields;

- Can provide help and explain different projects and problems in different areas;

- Can generate text - articles, program code, letters, poetry:

*1)* You can ask ChatGPT to write an article on any topic, specifying what tone or style to use - formal, casual, persuasive, descriptive, humorous, emotional, technical, and more.

*2)* Some programmers (IT students too) try to outsource the entire programming process to ChatGPT. It's not that this technology can't write good programming code, but it's still recommended to be used only as an additional tool in this area.

*3)* There is a free ChatGPT Writer extension for Gmail that can compose emails and messages by correcting grammatical errors, paraphrasing text, changing writing style, and summarizing text.

## III. PROBLEMS WHEN USING CHATGPT AND AI TECHNOLOGIES IN EDUCATION

The main group of problems in the use of AI technologies in education is ethical and, more precisely, problems related to plagiarism. If the lecturer and the students have the opportunity to use similar technologies, what will stimulate them to express their position and their opinions? These technologies provide a faster and easier way to create texts on a given topic or to solve set problems or tasks. The students sort of overdo their homework and in this way, they don't acquire the habit of writing and expressing their thoughts, they don't reason, and they don't look for an explanation for the problems they are given to solve, they don't put any thought into it. They may not even read the condition of the given problem or tasks, but simply use the copy-and-paste functions and get the result.

How to make students understand that it makes sense to know the definitions from the learning material so that they can search for information, respectively ask ChatGPT. It is clear that in the modern conditions of Internet access, it does not make sense for them to learn by heart and reproduce a text, it is important to be able to solve problems and tasks and, above all, to learn to think.

The collective opinion that ChatGPT will lead to the extinction of certain professions and thus put many people out of work is relatively popular and shared among the vast majority of people. There is a fear that it will replace the programmers, and more specifically - the junior programmers. However, if it does replace them because it solves elementary tasks perfectly, where will seniors come from if they have not been junior programmers? How will the seniors be so good if they've missed the moment of programming elementary tasks - while they were studying at the university they missed it, and then there was no way to work as such.

ChatGPT can find applications in the learning process as an intelligent assistant. Its particular advantage is that it can provide learners with interactive help at any time and from any place. The authors in [10] specified the following guidelines in which the use of ChatGPT may be useful to students:

- Provide information and resources, answer questions, organize information, help prepare for exams, and provide feedback;

- Improve language skills - grammar, vocabulary, and style during communication with the agent, as well as use ChatGPT to check their written text for syntactic and grammatical errors;

- Provides a new interactive way of learning languages - without restrictions on when and where, and has opportunities to generate realistic dialogues in a chosen and interesting area for the learner; can exercise their foreign language skills if they communicate with him in a chosen foreign language;

- Improving cooperation and communication - if students work in a team, the use of ChatGPT stimulates communication between the participants in the teams and also between them and the teachers;

- Provide support and motivation - ChatGPT can also act as a means of support and motivation for students. They can use ChatGPT to talk about their problems and concerns or ask for advice on how to better manage their time and tasks.

It should be noted that ChatGPT is not the only natural AI agent that can understand and generate conversation in natural human language. In February 2023, Google introduced Bard, which follows the LaMDA model but has similar features and applications [17, 18]. Our research is focused on ChatGPT and therefore does not describe other similar solutions.

## IV. RESEARCH METHODS

### A. Objectives and Contributions of the Experimental Research

The research aims to investigate the possibilities and extent of the use of ChatGPT by university students and final-year high-school students in the process of their education. Naturally, we consider all the risks and challenges of the unethical and illegal use of such tools in the learning process. The attitude of the students and their assessment of the capabilities of ChatGPT in the learning process are important because this would determine the use of these technologies in schools and universities. 102 surveyed students from the University of Ruse and students in the last year of Mathematical High School "Baba Ton-ka", Ruse, Bulgaria took part in the research and focused on the place and role of ChatGPT as a potential source of knowledge and information for students and students. The main questions that are the aim of the study are:

R1: How familiar are the students with the capabilities of ChatGPT?

R2: What are the potential benefits and challenges associated with using ChatGPT in learning from the learner's perspective?

R3: Are students inclined to use ChatGPT in the university/school and what do they think they will achieve by using it?

R4: Can learners rate the responses received from ChatGPT?

### B. Description of the Respondents

The total number of participants in the study is 102, students from the Computer Systems and Technologies specialty of the University of Ruse, Bulgaria, and students in the last year of Mathematical High School "Baba Tonka", Ruse, Bulgaria, who, in addition to mathematics, study informatics and information technologies intensively (see Table I). It is important to specify the major of the students and the subjects they are studying since it is very likely that their IT orientation has some certain influence on the way they accept these new technologies, as well as their natural greater interest in them.

This group of students is chosen because they are students, we have direct observations on. These are students who we teach, thus allowing us to get to know them better.

TABLE I. DEMOGRAPHIC STATISTICS

| Demographic Stats | | | |
|---|---|---|---|
| | *Possible Options* | *Number* | *Percentage* |
| Gender | Female | 44 | 43% |
| | Male | 56 | 55% |
| | I don't want to share | 2 | 2% |
| Age | Under 19 years (primarily high-school students) | 44 | 43% |
| | 19-25 years | 36 | 35.3% |
| | Over 25 years | 22 | 21.6% |

## C. Data Collection

An anonymous survey was developed using Google Forms. It is distributed among students and pupils through email messages, social networks, and messages on online learning platforms used in the university and school.

The survey consists of 15 questions of different types.

The questions can be divided into groups according to the information we expect from them. The first group of questions concerned various demographics of the respondents and their possible knowledge and experience in using AI and ChatGPT. Another group of questions concerns the respondents' assessment of the benefits and harms of using such technologies, and the challenges of using them, and also a group of questions to assess the extent to which students can judge how true the answers to the chat GPT are. Questions with different types of answers were used - an open answer that requires entering the opinion of the respondent in a free text, multiple choice of one of all the indicated answers, and choice of several of many possible answers.

## V. RESULTS

The first group of questions aims to answer the first research question - to assess the extent of knowledge and use of ChatGPT by the respondents. The results show that a very small percentage (13.7%) have not used ChatGPT at all, an equal number do not use it but have tried it and know what it is about. 21.6% regularly use it (see Fig. 2). In reality, only 13.7% of those polled don't use and are not interested in ChatGPT that shows that the technology is familiar to the trainees.

It is interesting what the respondents mostly use this technology for according to the answers of the respondents to the question "What kind of problem did ChatGPT help you with?", it was most used for research, during the development of projects in second place, and in third place for writing homework. An interesting fact from these answers is that it is also used during exams and 19% of the respondents state that it has helped them in some form during exams tasks and 17% of them - with answering exam questions. About a quarter still used ChatGPT to fill in ambiguities and gaps in their knowledge of the material they were studying.

Just over ⅓ of respondents rate ChatGPT as better than Google and other popular search engines, with almost half of them rating it not just better, but much better. Also, almost ⅓ don't care which is better, but rather getting an answer to their problem is important to them. A rather large percentage - 26.5% claim that they do not have any definite opinion on the

subject, and for 14.7%, ChatGPT is in no way superior to search engines (see Fig. 3).



Fig. 2. Answers to "Have you used ChatGPT in your school/university/work?"



Fig. 3. Answers to "Do you think ChatGPT is better than Google and other search engines?"



Fig. 4. Word cloud of users' opinions.

With this survey, we tried to find out what the respondents think about AI and ChatGPT and whether they think there is a difference between the terms by asking them to write in a few words their opinions. We summarized and evaluated the most frequently used phrases and words in the descriptions, and as can be seen in Fig. 4, they are aware of the differences between AI and ChatGPT, evaluating ChatGPT as a platform that uses AI. Most often rated ChatGPT as extremely useful, which they consider being the future of the field.

Regarding the second set of questions about the assessment of R3: Are students inclined to use ChatGPT in the university/school and what do they think they will achieve by using it?

As apparent from R1 the respondents surveyed have tried and used ChatGPT, it's not entirely new to them and some even have quite a bit of experience with it. We were interested to know if they were worried about having used the agent and what they think about the teaching knowing of their usage. A very large part of the respondents (53.9% do not see the point in sharing with the teacher the fact that they had to use this kind of help because they do not think it will improve the teaching or the content of the course and 15.7% are afraid to mention to the teacher, so as not to harm their final grade (see Fig. 5). In a small additional anonymous survey with ⅓ of respondents, regarding a specific homework assignment, 54.5% of participants admitted to using ChatGPT to write the source code implementing the assignment of the homework. 63.6% of the participants had to make corrections to the solution returned by ChatGPT, and the rest used it directly. 81.8% tested with different data and tried to fully verify the functionality of the returned code, and the rest admitted that they had not tested at all or attempted to test with any data, but rather trusted ChatGPT. These results show that trainees are coping and benefiting from using ChatGPT. It is worrying, however, that there are a considerable 20% who use it without thinking about the tasks set and the answers returned, and directly use them to pass them as a solution to homework.



Fig. 5. Answers "If you used ChatGPT in any discipline, did you tell the teacher?"

Another question in the survey gives an example of a small programming task that usually could be given as homework or in a workshop. The question itself states: "You have a task for finding the shortest path in a graph in C++…… Write the question that you would ask ChatGPT". The answers vary from complete copy-paste of the text of the task to breaking down the steps they would take to explain their problem to ChatGPT and to paraphrase the question. Some of the most interesting answers are:

- "No question just ctrl + c and ctrl + v on the task, He understands it and if my ideas are different from his answer then I ask specific questions, but this is rarely necessary."

- "I would like to give me the whole code so I can get the idea of how the algorithm works, then I'd try to write the code myself and if something doesn't work out I'd check from the already given answer."

- "Make me a road map for all the algorithms for finding the shortest path in a graph, ranking them for me from good and used to not so much, taking into account speed, complexity, and all advantages."

- "Write me some C++ algorithms that can find the shortest path in a graph from vertex A to vertex B."

- "Algorithms for shortest path in a graph C++".

It is interesting to note that from 102 answers, only 13 are in English. That gives the impression that the students are more willing to ask in their native language, in this case - Bulgarian, and receive the answers in the same language. While the ability of ChatGPT to understand and perceive different languages there are some cases where the agent is confused and thus susceptible to wrong answers. Fig. 6 shows a conversation with ChatGPT where the Bulgarian language is incorrectly recognized as Russian and the chat needs several interactions to understand and correct the problem.



Fig. 6. Screenshot of confused ChatGPT.

It is interesting to see if they need any additional guidance and pointers on how to use it if they are officially allowed to use AI agents. Only 8.8% felt they needed help at least in the beginning, and 66.7% would use it directly themselves without a problem (see Fig. 7).

Apart from whether they have used or would use it is interesting to see how they use ChatGPT and what types of questions they ask. The reason behind this is that a large proportion of the respondents have expressed their desire for copy able text assignments in the past. Students, and more specifically, the students who have responded to the survey, do not like assignments that are spoken and explained aloud. Here, it should be noted that the assignments in question, are simple tasks, meant to be done in the timespan of a workshop and not homework assignments.

Fig. 7. Answers "If the teacher allows you to use ChatGPT in your discipline, would you use it?"



Fig. 8. Answers to "What type of questions do you ask ChatGPT?"

Pretty quickly it came as obvious to this attachment to the text format. Given an assignment with 10 tasks, some students produce the answers in less than five minutes, which under normal circumstances is impossible, as it is physically impossible to write a programming source code for such a short time. Upon close inspection of the provided answers, it is not difficult to notice that they are generated by AI, as some small details are not described in the assignments and are known by the lecturers. This whole situation is quite an obvious sign of ChatGPT usage.

Fig. 8 shows that according to the self-assessment of the respondents, 46.1% try to ask guiding questions to the agent to orientate themselves in the problem and the topic based on the answers received. Some of them (18.6%) try to break the assignment into separate problems, ask questions related to them, and then summarize and combine the information from the received answers. There are still 18.6%, who don't bother to think about the task they have and directly give it to ChatGPT to get the ready answer. There are also 10% of respondents who do not have much success using the chatbot because they cannot ask their questions in a way that it understands them correctly.

Regarding the received answers is shown in Fig. 9. 42.2% of the respondents used the chat just for pointers about their problem and the same percentage found mistakes in the answers that they corrected before submitting or using them. A relatively small percentage, but still a notable percentage – 8.8% directly used the received answers without any corrections. This raises the question about their ability and willingness to further check and dive into the problem. 6.9% answered that they had not used the solution provided by ChatGPT. While the percentage is small, it is interesting to find out why is– maybe the provided solution was not correct at all or was too complex to understand and implement, or they have just wanted to see and play around with the chat.

R4: Can learners rate the responses received from ChatGPT - According to the results shown in Fig. 10, 30.4% of respondents do not bother to check the answers received from the chatbot, with 24.5% trusting them completely. Almost 20% check each answer further, and 50% first consider whether the answer can be used and only then check it.



Fig. 9. Answers to "Did you directly use the solution / answer that ChatGPT returned?"



Fig. 10. Answers of "Using ChatGPT do you check the answer in other sources?"

To assess how the respondents think about the answers from the chatbot, a question in which two definitions are given for the same term (in this case it is a definition of a set). All participants are aware of the term and use it in the learning and programming process. They are given a definition of the concept that is popular in textbooks and a definition given by ChatGPT. Their task is to evaluate the two definitions. The received answers show that for the respondents it does not matter whether the definition will be strictly formulated (in this case, it is the second definition, the one from the textbooks) or

it will be in a more descriptive form and with few examples (the first, from ChatGPT) - still for 38.2% the answer returned by ChatGPT is clearer and more understandable and they are more likely to trust it. For 31.4% the strictly theoretical versions is more understandable, and for 15.7% of the respondents, both are equally clear (see Fig. 11). What's worse is that for 14.7%, neither of the two definitions or explanations is comprehensible, and as we said - this is certainly something familiar to them, or at least should be.



Fig. 11. Answers to "Which of the two definitions below is more understandable to you?"

## VI. CONCLUSIONS

The usage of ChatGPT and similar AI-based technologies is something that is going to be more and more common. As lecturers and teachers, it is up to us to be able to navigate and adapt to it. The students have this very tempting, interesting, and easy-to-use at first glance technology. It is expected that they will be tempted to use it, after all this is an easy way to pass an exam and receive an excellent grade. While the education system does have its flaws and rewards excellent grades, the usage of these types of systems can and should be used in more efficient and effective manners to help with developing the thinking and problem-solving skills of the students.

From the experiment described in this article, the following can be concluded:

- ChatGPT systematizes sources of information found on the Internet on a given topic and saves time and effort;

- It provides personalized feedback and assistance to students anytime and from anywhere where they have access to the Internet;

- There is a real danger that students will learn false, malicious, or biased information if they rely entirely on ChatGPT without verifying the authenticity of what is written. As the survey shows, they do not pay enough attention and accept the answers as true;

- There is a real danger of fraud in the preparation of academic texts, cheating and plagiarism;

- The answers from ChatGPT can be deceiving and if the students are trusting it blindly, as this survey has proven to be the case, this can lead to bigger problems in the

future. Many of today's students are going to develop the habit of copy-and-pasting their problems in such chats and are going to stop developing their critical thinking thus limiting their intellectual growth.

This research underscores the undeniable potential of ChatGPT in reshaping educational dynamics. However, it also emphasizes the critical need for responsible integration. The findings spotlight the importance of equipping students with the skills to discern between AI-generated content and authentic knowledge. Striking a balance between leveraging AI for efficiency and preserving the essence of intellectual growth remains imperative in the evolving landscape of education.

While ChatGPT demonstrates remarkable capabilities, the study accentuates the irreplaceable role of human guidance in education. The findings highlight that, despite AI's potential to enhance learning experiences, it should be viewed as a tool rather than a substitute for human educators. The emphasis is on developing strategies to effectively integrate AI while ensuring that students receive the mentorship and critical thinking skills essential for their development.

As ChatGPT and similar AI technologies become integral to the educational experience, there is a pressing need to incorporate ethical AI education. The study underscores the importance of guiding students in understanding the ethical implications of relying on AI tools. Educators are encouraged to incorporate discussions on responsible AI use, fostering a generation that not only embraces technological advancements but also critically evaluates their impact on learning and intellectual development.

## REFERENCES

[1] George, A. S., & George, A. H. (2023). A review of ChatGPT AI's impact on several business sectors. Partners Universal International Innovation Journal, 1(1), 9-23.

[2] Kalla, D., & Smith, N. (2023). Study and Analysis of Chat GPT and its Impact on Different Fields of Study. International Journal of Innovative Science and Research Technology, 8(3).

[3] Felten, E., Raj, M., & Seamans, R. (2023). How will Language Modelers like ChatGPT Affect Occupations and Industries? arXiv preprint arXiv:2303.01157.

[4] Schön, E. M., Neumann, M., Hofmann-Stölting, C., Baeza-Yates, R., & Rauschenberger, M. (2023). How are AI assistants changing higher education? Frontiers in Computer Science, 5.

[5] Sallam, M., Salim, N., Barakat, M., & Al-Tammemi, A. (2023). ChatGPT applications in medical, dental, pharmacy, and public health education: A descriptive study highlighting the advantages and limitations. Narra J, 3(1), e103-e103.

[6] Chiesa-Estomba, C. M., Lechien, J. R., Vaira, L. A., Brunet, A., Cammaroto, G., Mayo-Yanez, M., ... & Saga-Gutierrez, C. (2023). Exploring the potential of Chat-GPT as a supportive tool for sial endoscopy clinical decision making and patient information support. European Archives of Oto-Rhino-Laryngology, 1-6.

[7] Biswas, S. (2023). Importance of chat GPT in Agriculture: According to chat GPT. Available at SSRN 4405391.

[8] Ausat, A. M. A., Suherlan, S., & Azzaakiyyah, H. K. (2023). Is ChatGPT Dangerous for Lecturer Profession? An In-depth Analysis. Jurnal Pendidikan Dan Konseling (JPDK), 5(2), 3226-3229.

[9] Ali, J. K. M., Shamsan, M. A. A., Hezam, T. A., & Mohammed, A. A. (2023). Impact of ChatGPT on learning motivation: teachers and students' voices. Journal of English Studies in Arabia Felix, 2(1), 41-49.

[10] Fauzi, F., Tuhuteru, L., Sampe, F., Ausat, A. M. A., & Hatta, H. R. (2023). Analysing the role of ChatGPT in improving student

productivity in higher education. Journal on Education, 5(4), 14886-14891.

[11] Shoufan, A. (2023). Exploring Students' Perceptions of CHATGPT: Thematic Analysis and Follow-Up Survey. IEEE Access.

[12] Rosenblatt, K. (2023). ChatGPT banned from New York City public schools' devices and networks. NBC News.

[13] Farrokhnia, M., Banihashem, S. K., Noroozi, O., & Wals, A. (2023). A SWOT analysis of ChatGPT: Implications for educational practice and research. Innovations in Education and Teaching International, 1-15.

[14] Lo, C. K. (2023). What is the impact of ChatGPT on education? A rapid review of the literature. Education Sciences, 13(4), 410.

[15] Neumann, M., Rauschenberger, M., & Schön, E. M. (2023). "We Need To Talk About ChatGPT": The Future of AI and Higher Education.

[16] Ausat, A., Massang, B., Efendi, M., Nofirman, N., & Riady, Y. (2023). Can Chat GPT Replace the Role of the Teacher in the Classroom: A Fundamental Analysis. Journal on Education, 5(4), 16100-16106. https://doi.org/10.31004/joe.v5i4.2745.

[17] Rayadurgam, V. C., & Afshan, N. Excuse Me ChatGPT & Bard! Can I Trust Your Data Analytical Skills?–a Academician vs Practitioner's Perspective. Can I Trust Your Data Analytical Skills.

[18] Destefanis, G., Bartolucci, S., & Ortu, M. (2023). A Preliminary Analysis on the Code Generation Capabilities of GPT-3.5 and Bard AI Models for Java Functions. arXiv preprint arXiv:2305.09402.

# From Time Series to Images: Revolutionizing Stock Market Predictions with Convolutional Deep Neural Networks

TATANE Khalid[1], SAHIB Mohamed Rida[2], ZAKI Taher[3]

National School of Applied Sciences, ESTIDMA, Ibn Zohr University, Agadir, Morocco[1]
Faculty of Applied Sciences, IMIS Laboratory, Ibn Zohr University, Agadir, Morocco[2, 3]

*Abstract*—**Predicting the trend of stock prices is a hard task due to numerous factors and prerequisites that can affect price movement in a specific direction. Various strategies have been proposed to extract relevant features of stock data, which is crucial for this domain. Due to its powerful data processing capabilities, deep learning has demonstrated remarkable results in the financial field among modern tools. This research suggests a convolutional deep neural network model that utilizes a 2D-CNN to process and classify images. The process for creating images involves transforming the top technical indicators from a financial time series, each calculated for 21 different day periods, to create images of specific sizes. The images are labeled Sell, Hold, or Buy based on the original trading data. Compared to the Long Short Time Memory Model and to the one-dimensional Convolutional Neural Network and the model exhibits the best performance.**

*Keywords—Technical indicators; convolutional neural networks; stock trend forecasting; deep learning*

## I. INTRODUCTION

Predicting stock trends has been the subject of intensive research for decades, making it a challenging problem [1] [2]. Developing accurate prediction models is challenging due to the many factors that can impact price changes. In recent years, machine learning and text mining models have shown promising results in multiple applications, including stock market prediction [3]. However, deep learning techniques [4] [5] [6] have become a potential alternative due to the limited accuracy of these models. The use of deep learning in financial studies is becoming increasingly popular as its potential is demonstrated in various applications. Despite this, deep neural networks are still not widely used for stock market prediction. The main reason for this is that processing financial data is difficult and the prediction task is complex. To tackle this issue, a unique strategy is suggested that transforms financial data into image-like representations and utilizes a deep convolutional neural network (CNN) to forecast stock trends.

This paper provides a thorough explanation of the approach, encompassing the technical aspects of the model and the data pre-processing procedures. The proposal consists of converting financial data into 2D images and then feeding them into a CNN for training. In particular, 21 technical indicators for a 21-day period are computed and the 225 most pertinent features are selected, which are transformed into 15x15

images. A chronological sequence of these images is used as input to the CNN model.

The approach's effectiveness in predicting stock trends is demonstrated by the presentation of its results. The model's performance is evaluated through various metrics, such as precision, confusion matrix, F1-score, and recall. Moreover, the evaluation of the outcomes of the proposed approach against other conventional machine learning models, proving that the deep learning approach is superior. In summary, the proposed approach is a promising solution for stock market prediction, delivering a more precise and efficient alternative to traditional machine learning models.

The rest of the paper is structured as follows: Section II describes the related work; Section III details the background and the followed methodology in Section IV. Section V covers the results, and Section VI presents general conclusions. Finally future works is given in Section VII.

## II. RELATED WORK

Different approaches have been investigated in the field of predicting stock prices to improve the accuracy of forecasts. [7] It suggests employing ARIMA models, for crop price predictions. It recommended using feature selection techniques to enhance their predictive abilities. The research in [8] developed an HNN-based forecasting approach that utilized headline information as well as ARIMA model predictions to drive sales estimation and business strategic decision making for the publishing industry.

For forecasting stock price indices, [9] came up with a method integrating ARIMA model and backward propagation neural network. The resulting statistics were compared against outputs attained via back propagation neural network and a single ARIMA model approach, respectively. Similar to [10], a model was designed to measure time periods optimal for the prediction of the market price of securities within various industries in 2019. The results emphasized the importance of using large datasets while training models to predict with more precision. A hybrid CNN-TLSTM model for predicting the USD/CNF exchange rate by [11] obtained low mean absolute percentage error (MAPE) and mean squared error (MSE) showing its ability of dealing with complex nonlinear data According to [12], CNN-LSTM proved efficient as a quantitative trading strategy analyzer in the stock market capable of forecasting future movement of stocks and revealing

their patterns. The study in [13] proposed the CNN-BiLSTM-AM technique for prediction of the following days' closing stock prices.

The study revealed that this model achieved greater predictive accuracy as well as performance when compared to other models using CNNs coupled with an attention mechanism. A study by [14] used an LSTM recurrent neural network to predict stock market trends, highlighting the LSTM model's capacity for accurately predicting stock prices.

## III. BACKGROUND

*A. Background*

Deep learning has become one of the most relevant mechanisms in the domain of machine learning and has been used generally in image recognition and natural language processing [15]. The most commonly used deep learning algorithms are convolutional neural networks (CNN) for feature extraction [16], recurrent neural networks (RNN) for mining time series data [17], etc. Therefore, for the purpose of making full use of their benefits, a CNN model is utilized to benefit from its power in feature extraction, which will be presented as follows.

*1) Convolutional neural networks:* Undoubtedly, convolutional neural networks (CNNs) (see Fig. 1) have established themselves as the most prevalent deep learning model, showcasing various network structures encompassing 1D, 2D, and 3D CNNs [18]. The 1D CNN specializes in processing sequential data [19], while the 2D counterpart excels in image and text recognition. Conversely, the 3D CNN takes the next stage in the identification of video data [20]. A standard convolutional neural network comprises an input layer, an output layer, and multiple hidden layers. These concealed layers typically house numerous convolutional layers that execute convolution through a dot product, employing RELU as the activation function. Following these

convolutional layers are additional operations such as pooling, flattening, and fully connected layers [21].

*2) Convolutional layer:* The convolutional rlayer, a cornerstone of convolutional neural networks, shouldes the bulk of computational burden. It utilizes a set of spatially-restricted filters, each spanning the entirety of the input volume's depth. During the forward pass, these filters are slid across the input's width and height, computing the dot product between their entries and the input at each position. This process, akin to fingerprint identification, generates a 2D activation map, revealing the filter's response at every location. The network learns filters that activate to specific visual features, each producing a distinct map. These maps are then stacked along the depth dimension, forming the layer's output volume.

*3) Pooling layer:* Convolutional neural networks frequently employ pooling layers interspersed between convolutional layers. These interludes serve to shrink the representation's size, thereby reducing both the network's parameter count and computational demands, while simultaneously mitigating overfitting. The pooling layer independently processes each depth slice of the input, resizing it through a MAX operation. A widely used variant employs 2x2 filters with a stride of two, effectively down sampling each depth slice by a factor of two in both width and height dimensions. This translates to each MAX operation extracting the highest value from a quartet of numbers, leaving the depth dimension unaltered.

*4) Fully connected layer:* Fully connected layers establish connections between every neuron in one layer and every neuron in another layer. This architecture is similar to that of a traditional multi-layer perceptron neural network (MLP). In a CNN, the flattened matrix is passed through a fully connected layer to enable classification.



Fig. 1. CNN model architecture.

## IV. METHODOLOGY

The proposal involves a model that employs Convolutional Neural Networks (CNNs) to detect favorable buying and selling points in stock prices. The model utilizes the top 15 technical indicators out of 21, using varying time intervals to generate images. As mentioned in (see Fig. 4), this approach encompasses four key stages: feature engineering, feature selection, data labeling, class imbalance handling, and image creation. The objective is to identify the most appropriate buy, sell, and hold positions in the time series of related stock prices.

*1) Feature engineering:* The dataset comprises features such as Date, Open, High, Low, Close, Adj Close, and Volume. To enhance the dataset's informational value, it's necessary to calculate 21 technical indicators for each day. These indicators include: RSI, William %R, MFI, MACD, PPO, ROC, CMFI, CMO, SMA, EMA, WMA, HMA, Triple EMA, CCI, DPO, KST, EOM, IBR, DMI, PSAR, and Volume values for intervals ranging from 6 to 27 days. These indicators are primarily categorized as momentum or oscillator types.

*2) Feature selection:* The proposed method improved the model's performance by applying a feature selection technique to the technical indicators. It selected 225 high-quality features using a univariate method that chooses variables with a strong relationship to the target outcome. It used two methods for feature selection: ANOVA F-value focusing on statistically significant differences between classes, and mutual information capturing non-linear relationships. It identified the features that both methods agreed on, then sorted the index list from the intersection of ANOVA F-value and mutual information to group similar indicators in the image.

*3) Labeling data:* The target was labeled by following a methodology [22] that used an 11-day window on the closing price. It checked if the middle number in the window was the highest, and if so, labeled the last day (11th day) as SELL. Otherwise, if the middle number was the lowest, labeled the last day as BUY. If neither condition was satisfied, labeled the last day as HOLD. The window was then slid and the process was repeated. This approach aimed to identify buying and selling opportunities at low and high points within each 11-day interval.

*4) Handling class imbalance:* The target was labeled and the data was found to be severely imbalanced. The Hold instances outnumbered the Buy and Sell instances by a large margin. The labeling algorithm used produced a relatively high number of Buy and Sell instances. Moreover, the Hold class was difficult to classify, making it hard for the model to learn meaningful insights. To solve this problem, the Sample Weights technique was adopted. This technique involved giving more weight to specific samples, which was a successful method for dealing with class imbalance.

*5) Image creation:* After completing all the preparatory steps, including downloading the dataset, calculating the technical indicators, performing feature selection, labeling the target, and normalizing the data, tabular data containing 225 features for each day was obtained. Subsequently, this data was converted into an image-like format. Sample 15x15 pixel images generated during this image creation process are presented in (see Fig. 2) to illustrate the visual representation of the data.



Fig. 2. Sample of generated images.

*6) CNN model architerture:* The proposed model consists of nine layers. The first layer, the input layer, receives images of size 15x15x3. The second layer is a convolutional layer with 25 filters, each of size 2x2. This layer is followed by a dropout layer that randomly drops out nodes during training with a 0.15% probability. Next, another convolutional layer is added, featuring 12 filters of the same size as the previous layer, stacked to a dropout layer with the same percentage (0.15%). Comes next the flatten layer, a dense layer with 100 neurons, a dropout layer with 0.15%, and the last layer which is the output layer with three neurons The model's output layer consists of three neurons, which determine the final action. The neural network uses the Adam optimizer, which is an adaptive gradient-based method that adjusts the learning rate for each parameter based on the previous gradients. It uses a learning rate of 0.001, which is a common value for the optimizer. The neural network trains for 3000 epochs, which means it sees the training data 3000 times with a batch size of 64 sample. EarlyStopping and ReduceLROnPlateau callbacks were utilized during model training. The first one stops the training when a monitored quantity (such as the validation loss) stops improving, and the last one reduces the learning rate when a monitored quantity (such as the validation loss) stops improving. (see Fig. 3) shows a recap for the CNN architecture.

Fig. 3.   Model.

*a) EarlyStopping:* Choosing the optimal number of training epochs is critical when training a neural network. Employing too many epochs can lead to overfitting, where the model becomes overly reliant on the training data and loses generalizability. Conversely, using too few epochs can result in underfitting, where the model fails to grasp the underlying patterns in the data. To address this dilemma, a technique called "early stopping" can be implemented. This involves training the model for a potentially large number of epochs, but proactively stopping the training process once the model's performance on a separate validation dataset ceases to improve.

*b) ReduceLROnPlateau:* During training, deep neural networks leverage the stochastic gradient descent (SGD) algorithm. This optimization technique utilizes randomly sampled data points from the training set to estimate the error gradient of the model's current state. Subsequently, the backpropagation algorithm is employed to update the model's weights based on this gradient information.

The learning rate governs how rapidly the model adapts to the problem it's facing. Smaller learning rates necessitate more training epochs because the weights are updated with smaller increments in each iteration. Conversely, larger learning rates propel faster adjustments, demanding fewer epochs.

However, selecting an inappropriate learning rate can lead to complications. A learning rate that's too high can cause the model to converge rapidly to a suboptimal solution, missing the optimal one. Conversely, a learning rate that's too low can cause the training process to become sluggish and potentially stall indefinitely. Consequently, selecting the optimal learning rate is crucial for effective deep neural network training.

*7) Performance evaluation:* The proposed model was assessed mainly by its computational performance, especially its accuracy in classifying buy, hold, and sell categories. These predicted labels guided buy, sell, and hold decisions based on stock prices. The model was trained and tested on WALMART stocks, using train, validation, and test splits. The F1-score was the main metric for performance evaluation on the test data, along with other metrics such as the confusion matrix, weighted F1 score.



Fig. 4.   Proposed method.

## V. RESULTS

On WALMART test data the model gave the following result:

The model's predictive abilities were thoroughly evaluated and analyzed. Table I presents the confusion matrix for the WALMART test data, while Table II offers a deeper analysis based on these results. Notably, the precision values for "Buy" and "Sell" classes surpass those for "Hold," indicating the model's strength in capturing key buy and sell points. However, false entry and exit points arose due to the inherent rarity of these instances compared to "Hold" periods. This posed a challenge for the neural network: accurately detecting rare entry and exit points while maintaining the overall "Hold" distribution. As a consequence, the model generated false alarms for non-existent entry and exit points to capture the majority of true ones. Additionally, the "Hold" points lacked the clear definition of "Buy" and "Sell" points, leading to some confusion, particularly near peaks and valleys within sliding windows.

TABLE I.    CONFUSION MATRIX OF TEST DATA (WALMART)

| Predicted | | | |
|---|---|---|---|
| Actual | Buy | Sell | Hold |
| Buy | 44 | 0 | 20 |
| Sell | 0 | 42 | 19 |
| Hold | 122 | 140 | 613 |

The F1-weighted score was used to comprehensively assess classification performance, which considered class weights for a more refined evaluation. The sample count of each class determined its F1 score weight. Multi-Layer Perceptron (MLP), Long Short-Term Memory (LSTM), and Convolutional Neural Network (CNN) methods were the baselines for stock market prediction, compared with the proposed model. Table III displayed the comparative prediction results.

TABLE II.    EVALUATION OF TEST DATA (WALMART)

| Total Accuracy: 0.699 | | | |
|---|---|---|---|
| | Buy | Sell | Hold |
| Recall | 0.6875 | 0.6885 | 0.7000 |
| Precision | 0.2650 | 0.2307 | 0.9401 |
| F1-Score | 0.3825 | 0.3455 | 0.8024 |
| Weighted F1 | 0.7480 | | |

TABLE III.    THE AVERAGE OF F1-SCORE OF TEST DATA (WALMART) ON DIFFERENT MODEL

| Model | AVG F1-Score |
|---|---|
| CNN | 0.44 |
| LSTM | 0.57 |
| MLP | 0.45 |
| **CNN 2D** | **0.74** |

## VI. CONCLUSION

In this study, a 2D Deep Convolutional Neural Network model was developed utilizing financial data and technical indicators to predict optimal buy and sell points. The model analyzed financial time series data, converted it into 2D images, and identified "Buy," "Sell," and "Hold" signals for profitable trades. Using WALMART stock prices as the primary data source, the model outperformed other models with the same dataset.

While the model's performance is encouraging, there's still room for improvement. Experimenting with various techniques, such as hyperparameter tuning, exploring different image sizes, alternative target labeling strategies, and utilizing additional technical indicators, could yield further gains. This research also contributes to the growing trend of applying 2D CNNs to non-image data, particularly in time series forecasting.

## VII. FUTURE WORKS

The model was trained on the historical data of WALMART stocks. For future work, more stock data from different sources will be incorporated to enhance the deep learning models. Moreover, other types of deep neural networks and their combinations will be explored to achieve better performance. Furthermore, since the news has an impact on the stock price trends, the model will be augmented with sentiment analysis features to increase its accuracy. Additionally, alternative labeling methods will be experimented to improve the model outcome.

## REFERENCES

[1] Kumbure, Mahinda Mailagaha, Christoph Lohrmann, Pasi Luukka, and Jari Porras. (2022). "Machine Learning Techniques and Data for Stock Market Forecasting: A Literature Review." Expert Systems with Applications 197: 116659.

[2] Guo, Y., He, F., Liang, C., & Ma, F. (2022). Oil price volatility predictability: new evidence from a scaled PCA approach. Energy Economics, 105, 105714.

[3] Henríquez, J., & Kristjanpoller, W. (2019). A combined Independent Component Analysis–Neural Network model for forecasting exchange rate variation. Applied Soft Computing, 83, 105654.

[4] Baek, Y., & Kim, H. Y. (2018). ModAugNet: A new forecasting framework for stock market index value with an overfitting prevention LSTM module and a prediction LSTM module. Expert Systems with Applications, 113, 457-480.

[5] Bao, W., Yue, J., & Rao, Y. (2017). A deep learning framework for financial time series using stacked autoencoders and long-short term memory. PloS one, 12(7), e0180944.

[6] Ma, C., & Yan, S. (2022). Deep learning in the Chinese stock market: The role of technical indicators. Finance Research Letters, 49, 103025.

[7] Shao, Y. E., & Dai, J.-T. (2018). Integrated Feature Selection of ARIMA with Computational Intelligence Approaches for Food Crop Price Prediction. 2018, 1-17. https://doi.org/10.1155/2018/1910520.

[8] Omar, H., Hoang, V. H., & Liu, D. R. (2016). A hybrid neural network model for sales forecasting based on ARIMA and search popularity of article titles. Computational Intelligence and Neuroscience, 2016. https://doi.org/10.1155/2016/9656453

[9] Du, Y. (2018). Application and analysis of forecasting stock price index based on combination of ARIMA model and BP neural network. In Proceedings of the 30th Chinese control and decision conference, CCDC 2018 (pp. 2854–2857). https://doi.org/10.1109/CCDC.2018.8407611

[10] Ghosh, A., Bose, S., Maji, G., Debnath, N. C., & Sen, S. (2019). Stock price prediction using lstm on Indian share market. EPiC Series in Computing, 63, 101–110. https://doi.org/10.29007/qgcz

[11] Wang, J., Wang, X., Li, J., & Wang, H. (2021). A prediction model of CNN-TLSTM for USD/CNY exchange rate prediction. IEEE Access, 9, 73346–73354. https://doi.org/10.1109/ACCESS.2021.3080459

[12] Liu, S., Zhang, C., & Ma, J. (2017). CNN-LSTM neural network model for quantitative strategy analysis in stock markets. Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), 10635(LNCS), 198–206. https://doi.org/10.1007/978-3-319-70096-0_21

[13] Lu, W., Li, J., Li, Y., Sun, A., & Wang, J. (2020). A CNN-LSTM-based model to forecast stock prices. Complexity, 2020. https://doi.org/10.1155/2020/6622927 Lussange, J., Palminteri, S., Bourgeois-Gironde, S., & Gutkin, B. (2019). Mesoscale impact of trader psychology on stock markets: A multi-agent AI approach.

[14] Moghar, A., & Hamiche, M. (2020). Stock market prediction using LSTM recurrent neural network. Procedia Computer Science, 170, 1168–1173. https://doi.org/10.1016/j. procs.2020.03.049

[15] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," nature, vol. 521, no. 7553, pp. 436–444, 2015

[16] Amziane, A., Losson, O., Mathon, B., & Macaire, L. (2023). MSFA-Net: A convolutional neural network based on multispectral filter arrays for texture feature extraction. Pattern Recognition Letters, 168, 93-99.

[17] Ullah, F., Bilal, M., & Yoon, S. K. (2023). Intelligent time-series forecasting framework for non-linear dynamic workload and resource prediction in cloud. Computer Networks, 225, 109653.

[18] J. Zhao, X. Mao, and L. Chen, "Speech emotion recognition using deep 1d & 2d cnn lstm networks," Biomedical Signal Processing and Control, vol. 47, pp. 312–323, 2019

[19] O. Abdeljaber, O. Avci, S. Kiranyaz, M. Gabbouj, and D. J. Inman, "Real-time vibration-based structural damage detection using one-dimensional convolutional neural networks," Journal of Sound and Vibration, vol. 388, pp. 154–170, 2017.

[20] H.-C. Shin, H. R. Roth, M. Gao, L. Lu, Z. Xu, I. Nogues, J. Yao, D. Mollura, and R. M. Summers, "Deep convolutional neural networks for computer-aided detection: Cnn architectures, dataset characteristics and transfer learning," IEEE transactions on medical imaging, vol. 35, no. 5, pp. 1285–1298, 2016.

[21] Wikipedia contributors, "Convolutional neural network — Wikipedia, the free encyclopedia," https://en.wikipedia.org/w/index.php?title=Convolutional_neural_network&oldid=962591055, 2023.

[22] [Sezer, O.B., Ozbayoglu, A.M., 2018. Algorithmic financial trading with deep convolutional neural networks: Time series to image conversion approach. Appl. Soft Computing. 70, 525-538.

APPENDIX

*1) Relative Strength Index (RSI):* The Relative Strength Index (RSI), conceived by J.W. Wilder in 1978, quantifies the relative intensity of upward price movements against downward ones. This technical indicator, employed by traders to assess trend strength and identify potential reversals, offers valuable insights into momentum and potential trading signals. As price fluctuations occur, RSI values oscillate between 0% and 100%, signifying whether prices reside in the overbought (above 70%) or oversold (below 30%) territory. Equation for calculating the RSI value is provided in Eq. (1):

$$RSI = 100 - \frac{100}{1+(averagegain/averageloss)} \qquad (1)$$

*2) Williams %R:* The Williams %R indicator, a stalwart among momentum-based technical tools, assists traders in identifying overbought and oversold conditions in stock prices. Its domain spans a spectrum from -100 to 0, serving as a window into market sentiment. A well-established interpretation of Williams %R values posits those readings below -80 signal an "oversold" market, ripe for potential upward reversals. Conversely, values

exceeding -20 suggest an "overbought" market, susceptible to downward corrections. Eq. (2) details the calculation of Williams %R.

$$R = \frac{\max(high)-close}{\max(high)-\min(low)} * -100 \qquad (2)$$

*3) Money Flow Index (MFI):* The Money Flow Index (MFI), a technical analysis stalwart, leverages both price and volume data to gauge overbought and oversold conditions in an asset, pinpointing potential trend reversals. Its values range from 0 to 100, offering a window into market sentiment. Unlike the Relative Strength Index (RSI), which focuses solely on price, the MFI encompasses both price and volume information, painting a more complete picture of market sentiment. Eq. (3) to Eq. (5) show the calculation of MFI:

$$RMI = TP \times Volume \qquad (3)$$

$$MFR = \frac{14\ PeriodPositiveMoneyFlow}{14\ PeriodNegativeMoneyFlow} \qquad (4)$$

$$MFI = 100 - \frac{100}{1+MFR} \qquad (5)$$

*4) Moving Average Convergence and Divergence (MACD):* The Moving Average Convergence Divergence (MACD) indicator is a technical analysis tool that gauges the momentum and direction of stock prices. It measures the relationship between two moving averages of a security's price to identify potential turning points in the market.

When the MACD line crosses above the signal line in an upward direction, it suggests that bullish momentum is increasing, and stock prices are likely to rise. Conversely, if the MACD line crosses below the signal line in a downward direction, it indicates that bearish momentum is gaining strength, and stock prices are likely to fall.

The MACD indicator is a versatile tool that can be used to identify potential buy and sell signals, assess the strength of a trend, and anticipate trend reversals. However, it's important to note that the MACD indicator should not be used in isolation and should be considered in conjunction with other technical indicators and fundamental analysis. Eq. (6) and Eq. (7) show the calculations of MACD and Signal Lines:

$$MACD_{Line} = EMA_{12} - EMA_{26} \qquad (6)$$

$$Signal_{Line} = EMA_9(MACD) \qquad (7)$$

*5) Percentage Price Oscillator (PPO):* While the Percentage Price Oscillator (PPO) might appear similar to the MACD, a closer look reveals a subtle difference. While both rely on moving averages to assess momentum and trend direction, their calculation methodologies diverge. Eq. (8) and Eq. (9) illustrate the specific computations employed by the PPO to arrive at its main line and signal line:

$$PPO = \frac{EMA_{12} - EMA_{26}}{EMA_{26}} \times 100 \qquad (8)$$

$$Signal_{Line} = EMA_9(PPO) \qquad (9)$$

*6) Rate of Change (ROC):* The Rate of Change (ROC) indicator shines a light on the velocity of price movements over a defined timeframe. Eq. (10) unveils its computational heart:

$$ROC = \frac{Lasted_{close} - Previous_{close}}{Previous_{close}} \times 100 \qquad (10)$$

*7) Chaikin Money Flow Indicator (CMFI):* The Chaikin Money Flow (CMF) indicator, a technical analysis stalwart, gauges the volume of capital flowing into or out of a security over time. Its values dance between 1 and -1, with proximity to 1 signifying robust buying pressure and proximity to -1 indicating potent selling pressure.

A CMF value hovering near 1 whispers of a substantial inflow of money into the security, suggesting buyers firmly grip the market reins. Conversely, a

CMF value teetering around -1 hints at a considerable outflow of capital, indicating sellers are asserting control.

However, it's crucial to acknowledge that the CMF indicator should not be a solitary oracle. Its effectiveness is amplified when combined with other technical indicators and fundamental analysis, creating a well-rounded investment perspective. For the mathematically inclined, the CMFs inner workings are unveiled in Eq. (11) to Eq. (13):

$$Multuplier = \frac{((Close-Low)-(High-Close))}{(High-Low)} \quad (11)$$

$$MFV(MoneyFlowVolume) = Volume \times Multiplier \quad (12)$$

$$CMF_{12period} = \frac{21periodSumOfMFV}{21periodSumOfVolume} \quad (13)$$

*8) Chande momentum oscillator (CMO):* The Chande Momentum Oscillator (CMO) joins the ranks of the RSI and the stochastic oscillator as a fellow momentum indicator. It too dances between -100 and 100, but with its own interpretive twist. When the CMO value waltzes past 50, stock prices are deemed to be frolicking in "overbought" territory. Conversely, a plunge below -50 suggests they're wallowing in the "oversold" zone. Eq. (14) unveils the CMO's inner workings:

$$CMO = 100 \times \frac{S_u - S_d}{S_u + S_d} \quad (14)$$

*9) Simple moving average (SMA):* The humble Simple Moving Average (SMA) reigns supreme as a fundamental indicator, capturing the essence of a security's price movement over a chosen period. Its power lies in its simplicity – it's just the average of prices during that timeframe. Its true magic unfolds when it's crossed by the actual stock price. Eq. (15) unveils the unassuming yet powerful formula behind the SMA:

$$SMA(M,n) = \sum_{k=a+1}^{a+n} \frac{M(k)}{n} \quad (15)$$

*10) Exponential moving average (EMA):* The exponential moving average (EMA) is a type of moving average indicator that places greater weight on recent price data, making it more responsive to current trends. Traders employ the EMA on trading charts to identify potential entry and exit points for trades based on the relative position of price action to the EMA. When the EMA is high relative to price, traders may consider selling, while when it is low relative to price, they may consider buying. Eq. (16) illustrates the calculation of EMA:

$$EMA(t) = (M(t) - x) \times \frac{2}{\tau+1} + x \quad (16)$$

$$\text{With } x = EMA(M, t-1, \tau) \quad (17)$$

*11) Weighted Moving Average (WMA):* While both the weighted moving average (WMA) and exponential moving average (EMA) aim to capture trends by assigning decreasing weights to past closing prices, their weighting mechanisms differ. The EMA employs an exponential decay, where the weights decline rapidly as we move back in time. In contrast, the WMA utilizes a linear decay, resulting in a more gradual decrease in weights assigned to older prices. This distinction highlights the potential for the WMA to offer a different perspective on market trends compared to the EMA, particularly in situations where recent price movements hold significant importance. Eq. (17) shows how WMA is calculated:

$$WMA(M,n) = \frac{SumOfWeighedAverages}{SumOfWeight} \quad (17)$$

*12) Hull Moving Average (HMA):* The Hull moving average (HMA) stands out among moving average indicators for its superior ability to minimize lag, a common drawback of traditional SMAs. While EMAs and

WMAs attempt to address this issue by emphasizing recent data, HMA takes it a step further, achieving even better lag reduction. Eq. (19) details the HMA's calculation, highlighting its distinct approach.:

$$HMA(M,n) = WMA\left(2 \times WMA\left(\frac{n}{2}\right) - WMA(n), \sqrt{n}\right) \quad (18)$$

*13) Triple Exponential Moving Average (TEMA):* The triple exponential moving average (TEMA) tackles a persistent challenge in technical analysis: how to filter out minor price fluctuations while capturing the underlying trend. Unlike traditional EMAs, TEMA employs a layered approach, essentially applying multiple exponential weightings to price data. This sophisticated strategy, as shown in the subsequent formula, effectively smooths out short-term volatility, revealing the true market direction with greater clarity:

$$TEMA = (3 \times EMA) - \left(3 \times EMA(EMA)\right) + \left(EMA\left(EMA(EMA)\right)\right) \quad (19)$$

*14) Commodity Channel Index (CCI):* Standing out as a versatile tool for traders, the Commodity Channel Index (CCI) assesses how much a current price deviates from its typical range over a specified period. Essentially, it compares the present price to the average price within that timeframe, translating the difference into an easily interpretable value. Although fluctuating between -100 and 100, the CCI spends a quarter of its time outside this range, signifying potentially significant deviations from normalcy. Eq. (20) and Eq. (21) unveil the mathematical magic behind this indicator, revealing how it transforms raw price data into a powerful market sentiment gauge:

$$CCI = \frac{TP-20PeriodsSMAofTP}{0.015 \times meanDeviation} \quad (20)$$

$$TypicalPrice(TP) = \frac{High+Low+Close}{3} \quad (21)$$

*15) Detrended Price Oscillator (DPO):* Unlike most oscillators chasing momentum, the Detrended Price Oscillator (DPO) takes a different path. Its unique strength lies in stripping away trend direction from price data, allowing it to focus solely on the underlying cyclical patterns. By filtering out trend noise, the DPO exposes the true rhythm of the market, highlighting peaks and troughs that hint at potential turning points. Eq. (22) unveils the secret behind this transformation, revealing how the DPO isolates cyclical swings for informed buy and sell decisions based on historical echoes:

$$DPO = \frac{PricefromX}{2} + 1periodAgo - XperiodSMA \quad (22)$$

Where X is the number of periods used for the look-back period.

*16) Know Sure Thing (KST):* he Know Sure Thing (KST) indicator, conceived by Martin Pring, addresses a fundamental challenge in technical analysis: deciphering the complexities of rate-of-change (ROC) data. KST tackles this challenge by employing a novel multi-period approach. It calculates simple moving averages (SMAs) of four distinct ROC periods, effectively smoothing out volatility and capturing momentum across various timeframes. These tamed ROC values then converge into a single, comprehensive indicator: the KST. To further enhance its interpretability, a 9-period SMA of the KST itself is employed as a signal line, generating crossover alerts for potential trading opportunities. Eq. (23) formalizes this innovative approach, unveiling the intricate interplay of SMAs and ROCs within the KST framework:

$$KST = RCMA_1 + (RCMA_2 \times 2) + (RCMA_3 \times 3) + (RCMA_4 \times 4) \quad (23)$$

where:

RCMA₁ = 10-period SMA of 10-period ROC

RCMA₂ = 10-period SMA of 15-period ROC

RCMA₃ = 10-period SMA of 20-period ROC

RCMA₄ = 15-period SMA of 30-period ROC

*17)Ease of Movement (EOM):* In the technical analyst's toolbox, the Ease of Movement (EOM) indicator stands out as a bridge between the ephemeral world of momentum and the tangible realm of volume. EOM seeks to capture this crucial intersection, condensing the combined influence of price movement and volume into a single, quantifiable value. This value, at its core, aims to answer a critical question: are prices moving "easily" in their current trend? EOM posits that effortless price movements, reflected in high EOM readings, hold the potential to persist and present lucrative trading opportunities. Eq. (24) to Eq. (26) unveil the mathematical alchemy behind EOM, revealing how it transforms raw price and volume data into this insightful blend of momentum and volume analysis:

$$DistanceMoved = (\frac{High+Low}{2} - \frac{PHigh+PLow}{2}) \qquad (24)$$

$$BoxRation = \frac{\frac{Volume}{scale}}{High-Low} \qquad (25)$$

$$1 - PeriodEOM = \frac{DistanceMoved}{BoxRation} \qquad (26)$$

*18)Interbank Rate (IBR):* The interbank rate, often overlooked yet crucial to the financial system's smooth functioning, operates like an invisible plumbing system connecting banks. This rate dictates the interest charged when banks borrow short-term funds from each other to meet immediate liquidity needs or conversely lend out excess reserves. These interbank loans, captured in Eq. (27), typically have very short durations, often overnight and rarely exceeding a week. Understanding the interbank rate is akin to comprehending the heartbeat of the financial system, as it regulates the flow of liquidity and influences borrowing costs for all institutions within the network:

$$IBR = \frac{SumOfInterestRates}{NbrOfInterestRates} \qquad (27)$$

*19)Directional Movement Indicator (DMI):* The Directional Movement Indicator (DMI) stands apart from conventional trend indicators by offering a multifaceted perspective. Instead of relying on a single metric, DMI deploys a three-pronged attack, harnessing the strengths of the Average Directional Index (ADX), the plus directional indicator (+DI), and the minus directional indicator (-DI). This powerful triumvirate delves deep into the trend's soul, simultaneously revealing its strength, direction, and potential turning points. DMI, with its values ranging neatly between 0 and 100, paints a clear picture of the market's underlying momentum, empowering traders to identify and capitalize on trending opportunities with greater confidence.

$$+DI = 100 \times EMA(\frac{+DMI}{ATR}) \qquad (28)$$

$$-DI = 100 \times EMA(\frac{-DMI}{ATR}) \qquad (29)$$

$$ADX = 100 \times EMA(\left|\frac{+DI - -DI}{+DI + -DI}\right| \qquad (30)$$

*20)Parabolic SAR:* The Parabolic SAR (SAR) acts as a vigilant sentinel, constantly scanning the market landscape for potential ambush points – where trends might stall or reverse. This dynamic indicator relies on a trio of key components: The previous SAR value, the extreme point (EP) marking the trend's peak or trough, and the ever-adapting acceleration factor (AF) reflecting the trend's sensitivity. As Eq. (31) and Eq. (32) illustrate, rising EPs fuel AF's ascent, amplifying the SAR's response to potential reversals. Conversely, falling EPs trigger AF's descent, making the SAR more cautious in a weakening trend. This intricate interplay between past, present, and future trends empowers traders to identify critical points of potential stops and reverses, navigating the market with newfound foresight.

$$PSAS + PresviousAF(PreviousEP + PSAR) \qquad (31)$$

$$PSAS - PreviousAF(PSAR - PreviousEP) \qquad (32)$$

# An Explainable and Optimized Network Intrusion Detection Model using Deep Learning

Haripriya C[1], Prabhudev Jagadeesh M.P[2]

Research Scholar, JSS Academy of Technical Education, Bengaluru, Affiliated to VTU Belagavi, India[1]
Assistant Professor, Global Academy of Technology, Bengaluru, Affiliated to VTU Belagavi, India[1]
Professor, JSS Academy of Technical Education, Bengaluru, Affiliated to VTU Belagavi, India[2]

*Abstract*—In the current age, internet and its usage have become a core part of human existence and with it we have developed technologies that seamlessly integrate with various phases of our day to day activities. The main challenge with most modern-day infrastructure is that the requirements pertaining to security are often an afterthought. Despite growing awareness, current solutions are still unable to completely protect computer networks and internet applications from the ever-evolving threat landscape. In the recent years, deep learning algorithms have proved to be very efficient in detecting network intrusions. However, it is exhausting, time-consuming, and computationally expensive to manually adjust the hyper parameters of deep learning models. Also, it is important to develop models that not only make accurate predictions but also help in understanding how the model is making those predictions. Thus, model explainability helps increase user's trust. The current research gap in the domain of Network Intrusion Detection is the absence of a holistic framework that incorporates both optimization and explainable methods. In this research article, a hybrid approach to hyper parameter optimization using hyperband is proposed. An overall accuracy of 98.58% is achieved by considering all the attack types of the CSE CIC 2018 dataset. The proposed hybrid framework enhances the performance of Network Intrusion Detection by choosing an optimized set of parameters and leverages explainable AI (XAI) methods such as Local Interpretable Model agnostic Explanations (LIME) and SHapely Additive exPlanations (SHAP) to understand model predictions.

*Keywords—Network Intrusion Detection; deep learning; hyper parameter optimization; hyperband; CSE CIC IDS 2018 dataset; XAI methods; LIME; SHAP*

## I. INTRODUCTION

With cyberattacks becoming increasingly prevalent, it is imperative that businesses shift their focus towards cybersecurity. Our lives have transitioned to be internet centric after the pandemic, but cybersecurity problems are also intensifying every day. Researchers are concentrating on creating Deep Learning (DL) based Network Intrusion Detection System (NIDS) to identify zero-day attacks as new varieties of cyberattacks continue to emerge. Outdated attack traffic and no representation of contemporary attack types leads to poor representation of real time network traffic. Also, redundancy, anonymity due to privacy or ethical issues, simulated traffic, lack of traffic diversity, and the absence of an all-inclusive dataset are some issues with most of the existing datasets. Despite numerous attempts, the research community is yet to accomplish the development of systems that can handle threats without human intervention. Malicious cyberattacks create significant security risks, necessitating the development of an innovative, adaptable, and more dependable Intrusion Detection System (IDS). The number of Internet-connected devices is anticipated to reach 50 billion by the end of the decade [1]. Although, techniques for infiltration and security defences have advanced dramatically during the past decade, a significant number of organizations still use outdated cybersecurity solutions.

Considering the above challenges, the primary objective of this research article, is to implement an explainable and optimized network intrusion detection model using DL techniques. The proposed work incorporates both optimization and XAI methods. The main contributions are as follows:

- Implement hyperband algorithm on the proposed DNNHXAI (Deep Neural Network Hypertuned XAI) model to choose optimized parameters.

- Investigate explainability of the proposed model using LIME and SHAP.

## II. LITERATURE SURVEY

Notable numbers of works are proposed in the area of NIDS, Abdulnaser et al. used Apache Spark and DL models on the CSE CIC 2018 dataset. The authors conclude that Spark drastically reduces the training time when compared to DL models [2]. Haripriya et al. performed distributed training of deep auto-encoder including all the attacks of the CSE CIC IDS 2018 dataset. The authors achieved an accuracy of 98.96% by training their proposed model on two worker nodes. [3]. Kanimozhi et al. used Artificial Neural Networks (ANN) and used only the benign and botnet traffic of the CSE CIC 2018 dataset. They used Grid Search CV to perform hyper parameter tuning. However, the authors conclude that their proposed model can be extended to detect the remaining classes of the dataset and usage of higher end frameworks like Tensor Flow to perform hyper tuning optimization [4]. Vimal Gaur implemented Machine Learning (ML) algorithms on CICDDoS2019 dataset to detect Distributed Denial of Service (DDoS) attacks and performed hyper parameter tuning. The author concludes that hyper parameter tuning increases the accuracy by 2.01% [5]. Priya Maidamwar et al. implemented ML algorithms on UNSW-NB15 dataset and used Grid Search CV as their hyper parameter tuning technique. An improvement in accuracy and minimization of False Alarm Rate (FAR) was observed [6].

Amin et al. propose a ML based NIDS model for binary classification on CSE CIC-IDS 2018 dataset. However, they observed that minority classes are misclassified due to class imbalance and suggest researchers to use techniques like Synthetic Minority Oversampling Technique (SMOTE) [7]. Haripriya et al. effectively addressed the class imbalance problem of the CIC CSE IDS2018 dataset by using SMOTE. They used deep autoencoder to classify all the attacks of the dataset [8]. Rambasnet applied and compared various State-of-the-art frameworks on CSE CICIDS2018 dataset. Their findings demonstrate the usefulness of different DL frameworks for detecting network intrusion traffic. An accuracy of 99% was achieved. However, since the class imbalance of the dataset was not addressed, a large number of infiltration samples were misclassified [9].

Anita Shiravani et al. proposed a new method for effectively selecting features using fuzzy numbers. The authors emphasize on the fact that dimensionality reduction plays a major role in pre-processing which in turn improves the system performance [10]. Mohammad Mausam et al. proposed a NIDS framework using Bayesian Optimization (BO) with Gaussian Process (GP). They implemented their proposed method on NSL-KDD dataset and conclude BO-GP outperforms Random Search Optimization [11]. Yoon Teck et al. implemented ML algorithms on CICIDS 2017 dataset and used BO-Tree-structured Parzen Estimator (BO-TPE) as the hyper parameter tuning technique. The authors direct future researchers to apply hyper parameter optimization on DL algorithms to substitute their ML approaches [12].

Abdulatif et al. implemented ML algorithms in Kitsune dataset used in the domain of NIDS. The authors recommend Grid Search optimizer with Tree algorithm for Kitsune dataset [13]. Hyojoon et al. use Proximal Policy Optimization (PPO) algorithm on CICDS2017 and UNSW-NB15 datasets to control the hyper parameters of Deep Neural Network (DNN)-based feature extractor and K-Means cluster module. The authors conclude that feature engineering is crucial in NIDS data pre-processing and direct future researchers to carry out research using diverse datasets [14]. Sara Emadi et al. implement Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) algorithms on the NSL-KDD dataset. The authors conclude that research in the area of NIDS can further be improved by reducing the training time and using hyper parameter tuning to improve the overall performance [15]. Haripriya et al. carry out a comprehensive study on different benchmark IDS datasets and their impact on network intrusion techniques [16]. The authors insist on the fact that the quality of the dataset plays a vital role in the domain of NIDS.

Zhibo zang et al. carries out an extensive survey on different methods, categorization, research gaps and challenges of XAI in the domain of Cyber security [17]. Pieter Barnard et al. use XGBoost model on the NSL-KDD dataset. The authors use SHAP to explain their proposed model [18]. Zakaria et al. use DNN on NSL-KDD and UNSW-NB15 dataset. The authors conclude their work by using LIME, SHAP and Rule fit methods to improve the interpretability of the proposed model [19]. Shraddha Mane et al. implement DNN on NSL-KDD dataset and use XAI methods to generate explanations [20]. Basim Mahbooba et al. addressed explainability by using Decision Tree (DT) on KDD dataset [21]. Syed wali et al. implements Random Forest (RF) on CIC CSE IDS2018 dataset and used SHAP as an XAI method [22].

Studies from literature reveal that hyper parameter tuning is very important to decide on the best architecture of the DL model. Although researchers have previously worked on optimization algorithms for hyper parameter tuning, there is an increasing need to use advanced optimization algorithms for hyper parameter tuning to speed up the training process. When it comes to optimization, overlooking hyperparameter optimization altogether is the most substantial mistake one can make. Modest adjustments to hyperparameter values can have a significant effect on model's performance. Especially in the domain of network security, the main aim is to speed up the process of intrusion detection and help network administrators to take immediate action before a catastrophic attack on the network occurs. XAI methods can help minimize model bias by outlining the standards for making decisions. Therefore, monitoring the model using XAI methods help in lessening the bias and also unexpected consequences. Previous research works in the field of NIDS reveal that, researchers have either used hyper parameter optimization methods or XAI methods in the field of NIDS. It is observed that none of the researchers used both the methods. This research article leverages both hyper parameter tuning and XAI methods, thus providing hyper optimization along with a comprehensive understanding of model's predictions.

## III. METHODOLOGY

The proposed framework is divided into two stages. First hyper parameter tuning is done using hyperband. Secondly, XAI methods such as SHAP and LIME are used to interpret the model predictions.

### A. Dataset Description and Pre-processing

The proposed work uses the CSE CIC IDS2018 dataset. The main rationale behind choosing the CSE CIC IDS 2018 dataset is, it reflects the current attacks. Unlike the outdated KDD Cup 99 dataset, it also includes a wide range of attacks. It is very essential to choose the dataset that reflects real time network traffic comprising a wide diversity of attacks. The dataset consists of a total of 16,000,000 samples spread over 10 CSV files. The dataset was collected from Amazon Web Services (AWS) S3 bucket [23].

Data type conversion was carried out by converting 64 bit values to 32 bit values. Features containing only one variable were dropped. Also features having infinite and Not a Number (NAN) values were also dropped. Label encoding and one-hot encoding was performed on the different attack types of the dataset. It was noticed that the dataset suffers from class imbalance. It is observed that the number of samples belonging to benign (normal) were more when compared to the attack class. Class imbalance leads to "Accuracy paradox". For instance, while using training data with a very high percentage of benign samples, a model could be trained to predict normal traffic with high accuracy, but it might not be good at detecting attack traffic. Similar observations were made on all the files of the dataset. Thus, to overcome the

limitations of the class imbalance, SMOTE was used in the proposed model.

### B. Workflow of the Proposed Model

The entire workflow of the proposed framework is shown in Fig. 1. The CSE CIC IDS 2018 dataset is collected from the AWS S3 bucket. To help speed up the training process and improve accuracy, pre-processing is carried. Class imbalance problem of the dataset is addressed using SMOTE [24]. Then the DL model is developed by using DNN. Hyper parameter tuning is carried out using hyperband algorithm and an optimized set of parameters are chosen [25]. The performance of the model is evaluated on the test dataset.



Fig. 1. Workflow of the proposed model.

### C. Importance of Hyper Parameter Tuning in the Proposed Work

The performance of any ML/DL model depends on the configuration. One major challenge in implementing any ML/DL algorithm is discovering an optimal configuration for the model and the training algorithm. Hyper Parameter Optimization (HPO) is a technique to deal with the challenge of fine-tuning DL hyper parameters. Tuning in an enormous search space is an exhausting process. Data-driven techniques must be employed to address the issues with HPO. Manual processes are ineffective. There are several hyper parameter tuning algorithms namely Random search, Bayesian Optimization (BO) and hyper band. Random search is the least efficient algorithm as it randomly selects parameter combinations from a search space rather than learning from previously tried parameter combinations. BO uses an optimization method that is sequential, and thereby it cannot be used well with parallel resources pair.

Speeding up configuration evaluation is the primary objective of an orthogonal approach to hyper parameter optimization. Hyperband can be considered as an extension of the successive halving approach; the goal of the Hyperband is to regularly apply successive halving to address the trade-off between the number of configurations and resource allocation. Additionally, it can find the ideal combination faster by using successive halving. The primary idea is to fit numerous models for a limited number of epochs and to only continue training the models that perform best on the validation set. Therefore, in comparison to commonly used hyper tuning algorithms like BO, hyperband can dramatically speed up a variety of DL and kernel-based learning tasks. All the above factors motivated us to use hyperband as the hyper parameter tuning technique in the proposed model.

### D. Need for Model Explainability in the Proposed Work

Considering the large sizes of NIDS datasets, performance becomes the bottleneck. DL models are incomprehensible, counterintuitive, and challenging for people to understand. All the DL models act like black-box structures. Because DL models are so complicated, interpretability research has taken multiple avenues. Over the years, DL models evolved by improving the performance metrics to handle large data but with increasing complexity came less interpretability. Feature importance methods were used to show how each feature is important to model prediction in general. However, these methods do not give information about individual predictions. Also, which features tend to increase or decrease the prediction is not known. Understanding ML model is referred to as model explainability. There are numerous advantages of integrating XAI methods with DL algorithms. It enables individuals to mitigate the negative impacts of automated decision-making and help in more informed decisions. To identify and protect security vulnerabilities. Integrating algorithms with human values is an essential goal. As an instance, suppose a model is able to determine if the traffic is normal or malicious, the network administrators have to know what parameters the model has considered. This helps to know whether the model contains any bias. It is also essential for network administrators to understand and describe the model's predictions once it has been implemented. Fig. 2 illustrates the importance of XAI methods in the proposed work.



Fig. 2. Explainable AI methods used in the proposed DNNHXAI model.

## IV. EXPERIMENTAL SETUP

Table I illustrates the different hyper parameters used while training the proposed DNNHXAI model. The number of hidden units ranged from 2 to 32 with a step value of 3. The number of hidden values ranged from 2 to 10. Different activations such as relu, tanh, sigmoid were used. Relu was the most preferred activation function. The dropout values ranged from 0.0 to 0.1 with a step size of 0.05. An optimized learning rate of 0.001 was chosen. Table II gives the different general parameter values used while training. The batch size was set to 128 with 15 epochs. The loss functions used were binary cross entropy and categorical cross entropy for binary and multi-classification respectively. Adam optimizer was chosen as preferred optimizer as it helps the model to converge faster.

TABLE I. DIFFERENT PARAMETERS USED FOR HYPER PARAMETER TUNING

| Sl No | Name of the hyper parameter | Range of values for different hyper parameters | Best hyper parameters given by Hyperband |
|---|---|---|---|
| 1 | Number of units | Min_value=2, Max_value = 32 Step=3, Default=32 | Units in $0^{th}$ layer = 29 Units in $1^{st}$ layer = 5 |
| 2 | Number of layers | Min_value =2, Max_value = 10 | 2 |
| 3 | Activation | Dense Activation Values=relu,tanh, sigmoid Default= relu | relu |
| 4 | Dropout | Min_value=0.0, Max_value = 0.1 Default = 0.005, Step = 0.05 | 0.1 |
| 5 | Learning Rate | Values = 1e -2,1e -3,1e -4 | 0.001 |

TABLE II. GIVES THE DIFFERENT GENERAL PARAMETER VALUES USED WHILE TRAINING

| Sl.No | Parameter | Value |
|---|---|---|
| 1 | Batch size | 128 |
| 2 | Number of epochs | 15 |
| 3 | Loss function | Binary Cross entropy Categorical Cross entropy |
| 4 | Optimizer | Adam |

All the experiments were carried out using Google Colab which is a cloud-based environment. To speed up the training process, Graphical Processing Units (GPU) was chosen as the runtime option. The train – test split was set to 75% and 25 % respectively.

## V. RESULTS AND DISCUSSION

An accuracy of 96.67% was achieved without hyper parameter tuning. With usage of hyperband (see Fig. 3) as the hyper parameter tuning technique, the accuracy peaked to 98.58%.

SHAP and LIME methods are used to explain the predictions of the proposed model [26] [27]. Fig. 4 gives the waterfall plot. It shows how a positive SHAP value positively impacts the prediction. On the contrary, a negative SHAP

Value has a negative impact on the prediction. The magnitude helps us understand how strong the impact is. It also illustrates the feature importance of SHAP analysis by using the summary plot by considering the CSV file containing Distributed Denial of Service (DDOS) attack. The chosen CSV file contains two classes benign (normal traffic) and DDOS (attack traffic). The Class label is encoded as '0' and '1' for Benign and DDOS attack respectively. As per the result, min packet length is the highest ranking feature.



Fig. 3. Difference in accuracy with and without hyperband optimization.

Fig. 5 illustrates how LIME can be used to understand local predictions given by the model by considering the Comma Separated Values (CSV) file containing DDOS attack in the CSE CIC IDS 2018 dataset. The features shaded in blue indicate positive influence on the output. Conversely, the features shaded in orange indicate negative influence on the output. Similar experiments were conducted on the different attacks of the dataset. The key difference between SHAP and LIME is how they provide explanations. SHAP uses a game-theoretic approach to provide global explanations. Conversely, LIME is model specific that provides local interpretable explanations. In this research work, an attempt was made to investigate model interpretability using SHAP and LIME. However, it is observed that LIME explanations are not robust because of its instability. For each prediction, a new explanatory is generated by the LIME algorithm. Thus, small variations in the data lead to different interpretations. In contrast, SHAP helps in providing global explanations, therefore explaining the overall model's behavior across all the instances. Finally, we conclude that SHAP performs better than LIME.

Table III gives the comparative analysis of the proposed work with other latest works exiting in the literature. It is observed that researchers have either used Optimization or Explainability but not both. Also, outdated datasets like NSL-KDD that do not reflect current attacks are still being used. Conventional hyper parameter tuning techniques like Grid Search CV are no longer suitable as it is time-exhaustive and computationally expensive, especially if it involves a high dimensional search space. Although, Random Search CV is better than Grid search CV, a lot of variance is observed because of its randomness. Research works [3][4][5][10][11] use different optimization methods for hyper parameter tuning. Research works [17][18][19][20][21][28] use different XAI methods. Table III clearly illustrates that none of the previous works in the field of NIDS incorporated a hybrid model leveraging both hyper-parameter optimization and explainability. Comparison was based on the usage of optimization, XAI method and accuracy as the performance

metric. In this research article, a hybrid approach incorporating both optimization and explainability is implemented. Advanced Optimization algorithms such as a hyperband helps in finding the hyper parameters faster with

improved accuracy. Explainable methods such as LIME and SHAP help in gaining greater insights on the data by understanding model predictions and thus increasing user's trust in the model.



Fig. 4. SHAP explanations using summary and waterfall plot.



Fig. 5. LIME local explanations.

TABLE III. COMPARATIVE ANALYSIS OF THE PROPOSED WORK

| Sl.No | Authors | Algorithm used | Dataset Used | Optimization method | XAI method | Accuracy |
|---|---|---|---|---|---|---|
| 1 | Kanimozhi et al. [3], 2019 | ANN | CSECIC IDS 2018 | Grid Search CV | - | 99.97% |
| 2 | Vimal Gaur et al. [4], 2022 | ML algorithms | CICDDoS2019 | NA | - | 98.78% |
| 3 | Priya R Maidamwar et al. [5] , 2022 | RF and MLP | UNSW NB15 | Grid Search CV | - | 99.34% |
| 4 | Mohammad Mausam et al. [10], 2022 | DNN | KDDTest+ KDDTest21 | BO-GP BO-GP | - | 82.95% 54.99% |
| 5 | Yoon Teck et al. [11], 2022 | ML algorithms | CICIDS 2017 | BO-TPE | - | 98% |
| 7 | Pieter Barnard et al. [17], 2022 | XGBoost, autoencoder | NSL -KDD | - | SHAP | 93.28% |
| 8 | Zakaria et al. [18], 2022 | DNN | NSL-KDD and UNSW-NB15 | - | LIME, SHAPE, and Rule Fit | 88% |
| 9 | Shraddha Mane et al. [19], 2021 | DNN | KDD test+ | - | SHAP, LIME, and BRCG | 82.4% |
| 10 | Basim Mahabooba [20], 2021 | DT | KDD | - | Self-explainable | NA |
| 11 | Syed Wali et al. [21] , 2021 | Stacked RF | CICIDS | - | SHAP | 98.5% 100% |
| 12 | Deepak Kumar et al. [28], 2022 | RF, KNN | NSL KDD99 | - | SHAP, LIME | 99.4% |
| **13** | **Proposed Work** | **DNNHXAI** | **CSECICIDS 2018** | **Hyperband** | **SHAP, LIME** | **98.58%** |

## VI. CONCLUSION

With the advancement in technology, the number of cyberattacks is increasing exponentially. Although, DL models prove to be efficiently detect intrusions, its complexity has increased tremendously at the price of massive computational overhead. It is exhausting, time-consuming, and computationally expensive to manually adjust the hyper parameters of DL models. In this research paper, hyperband an advanced hyper parameter tuning algorithm is applied on the proposed DNNHXAI model. It is observed that the configuration of model hyper parameters has a significant impact on its prediction accuracy. Although DL models today are able to achieve very good accuracies, there is an increasing need to enhance the user's trust by using XAI methods. First, the algorithm should have the best performing parameter configured and XAI methods should be used to deduce the contributing factors. Particularly, in the domain of cybersecurity, an attacker can largely exploit a vulnerability within few seconds. To address the above stated challenges, an attempt is made to not only configure the best parameters but also to understand the model predictions in an efficient manner. A single model that can detect a variety of attacks is proposed. It is efficient to quickly differentiate between normal and attack traffic. The proposed model overcomes the problems encountered in traditional DL algorithms w.r.t hyper parameter optimization and explainability. Instance by instance explanation is done with both LIME and SHAP. The main outcome of combining hyper parameter tuning with XAI techniques is to enable network administrator to take appropriate action based on the certainty of a detected attack. Considering all the files of the dataset, an overall accuracy of 96.67% and 98.56% is achieved without and with hyper parameter tuning respectively. The framework implements efficient pre-processing techniques, addresses class imbalance, uses the latest benchmark IDS dataset that reflects recent attacks, implements advanced hyper parameter tuning

techniques and leverages XAI methods to understand model's predictions. Promising results were achieved and an improvement in model's performance is observed when hyper parameter tuning is used. XAI methods are used to increase the explainability of model's predictions. As a future work, researchers are advised to leverage transfer learning techniques on the latest datasets in the domain of NIDS. Also, additional XAI methods can be used on different DL algorithms to explain model's predictions more efficiently.

## REFERENCES

[1] Singh, Satyanand. (2021). Environmental Energy Harvesting Techniques to Power Standalone IoT-Equipped Sensor and Its Application in 5G Communication. Emerging Science Journal. 4. 116-126. 10.28991/esj-2021-SP1-08.

[2] Hagar, Abdulnaser & Gawali, Dr.Bharti. (2022). Apache Spark and Deep Learning Models for High-Performance Network Intrusion Detection Using CSE-CIC-IDS2018. *Computational Intelligence and Neuroscience.* 2022. 1-11. 10.1155/2022/3131153.

[3] Haripriya C and Prabhudev Jagadeesh M. P, "Distributed Training of Deep Autoencoder for Network Intrusion Detection" *International Journal of Advanced Computer Science and Applications* (IJACSA), 14(6), 2023. http://dx.doi.org/10.14569/IJACSA.2023.0140633.

[4] Kanimozhi, V. & Jacob, Prem. (2019). Artificial Intelligence based Network Intrusion Detection with hyper-parameter optimization tuning on the realistic cyber dataset CSE-CIC-IDS2018 cloud computing. ICT Express. 5. 10.1016/j.icte.2019.03.003.

[5] Vimal Gaur, Dr. Rajneesh Kumar (2022). HPDDoS: A HyperParameter Model for Detection of Multiclass DDoS Attacks. *Vol. 71 No. 3s2 (2022): Special Issue on Mathematics Theory and its Contribution in Robotics and Computer Engineering.*

[6] Maidamwar, P. R., Bartere, M. M., & Lokulwar, P. P. (2022). Classification of Hybrid Intrusion Detection System Using Supervised Machine Learning with Hyper-Parameter Optimization. *Journal of Algebraic Statistics,* 13(3), 1532-1550.

[7] Amin Lama, Dr. Preeti Savant (2022). Network-Based Intrusion Detection Systems Using Machine Learning Algorithms. *International Journal of Engineering Applied Sciences and Technology,* 2022 Vol. 6, Issue 11, ISSN No. 2455-2143, Pages 145-155.

[8] Haripriya C, Prabhudev Jagadeesh M. P (2022). An Efficient Autoencoder Based Deep Learning Technique to Detect Network

Intrusions. *International Transaction Journal of Engineering, Management, & Applied Sciences & Technologies,* 13(7), 13A7P, 1-9. http://TUENGR.COM/V13/13A7P.pdf DOI: 10.14456/ITJEMAST.2022.142.

[9] Basnet, Ram & Shash, Riad & Johnson, Clayton & Walgren, Lucas & Doleck, Tenzin. (2019). Towards Detecting and Classifying Network Intrusion Traffic Using Deep Learning Frameworks. 10.22667/JISIS.2019.11.30.001.

[10] Shiravani, A., Sadreddini, M.H. & Nahook, H.N. (2023) Network intrusion detection using data dimensions reduction techniques. *J Big Data 10*, 27. https://doi.org/10.1186/s40537-023-00697-5.

[11] Masum, Mohammad & Shahriar, Hossain & Haddad, Hisham & Hossain Faruk, Md Jobair & Valero, Maria & Khan, Md & Rahman, Mohammad & Adnan, Muhaiminul & Cuzzocrea, Alfredo. (2022). Bayesian Hyperparameter Optimization for Deep Neural Network-Based Network Intrusion Detection..

[12] Yoon-Teck Bau, Tey Yee Yang Brandon. (2022) Machine Learning Approaches to Intrusion Detection System Using BO-TPE. Atlantis Highlights in Computer Sciences. *Proceedings of the International Conference on Computer, Information Technology and Intelligent Computing* (CITIC 2022).

[13] Alabdulatif, & Rizvi, Sajjad. (2023). Network intrusion detection system using an optimized machine learning algorithm. Mehran University Research Journal of Engineering and Technology.42.153.10.22581/muet1982.2301.14.

[14] Han, H.; Kim, H.; Kim, Y. An Efficient Hyperparameter Control Method for a Network Intrusion Detection System Based on Proximal Policy Optimization. Symmetry 2022, 14, 161. https://doi.org/10.3390/sym14010161.

[15] Al-Emadi, Sara & Al-Mohannadi, Aisha & Al-Senaid, Felwa. (2019). Using Deep Learning Techniques for Network Intrusion Detection. 10.1109/ICIoT48696.2020.9089524.

[16] Haripriya C, Prabhudev Jagadeesh M. P. "A Review of Benchmark Datasets and its Impact on Network Intrusion Detection Techniques," 2022 Fourth International Conference on Cognitive Computing and Information Processing (CCIP), Bengaluru, India, 2022, pp. 1-6, doi: 10.1109/CCIP57447.2022.10058660.

[17] Zhang, Zhibo & Al Hamadi, Hussam & Damiani, Ernesto & Yeun, Chan & Taher, Dr. Fatma. (2022). Explainable Artificial Intelligence Applications in Cyber Security: State-of-the-Art in Research. 10.48550/arXiv.2208.14937.

[18] Barnard, Pieter & Marchetti, Nicola & Silva, Luiz. (2022). Robust Network Intrusion Detection Through Explainable Artificial Intelligence (XAI). IEEE Networking Letters. 4. 1-1. 10.1109/LNET.2022.3186589.

[19] Z. A. E. Houda, B. Brik and L. Khoukhi, ""Why Should I Trust Your IDS?": An Explainable Deep Learning Framework for Intrusion Detection Systems in Internet of Things Networks," in IEEE Open Journal of the Communications Society, vol. 3, pp. 1164-1176, 2022, doi: 10.1109/OJCOMS.2022.3188750.

[20] S. Mane and D. Rao, "Explaining Network Intrusion Detection System Using Explainable AI Framework." arXiv, Mar. 12, 2021. doi: 10.48550/arXiv.2103.07110.

[21] B. Mahbooba, M. Timilsina, R. Sahal, and M. Serrano, "Explainable Artificial Intelligence (XAI) to Enhance Trust Management in Intrusion Detection Systems Using Decision Tree Model," *Complexity*, vol. 2021, p. e6634811, Jan. 2021, doi: 10.1155/2021/6634811.

[22] Wali, syed & Khan, Irfan. (2021). Explainable AI and Random Forest Based Reliable Intrusion Detection system. 10.36227/techrxiv.17169080.v1.

[23] A Realistic Cyber Defense Dataset (CSE-CIC-IDS2018) was accessed on 01.01.2023 from https://registry.opendata.aws/cse-cic-ids2018.

[24] Chawla, Nitesh & Bowyer, Kevin & Hall, Lawrence & Kegelmeyer, W. (2002). SMOTE: Synthetic Minority Over-Sampling Technique. J. Artif. Intell. Res. (JAIR). 16. 321-357. 10.1613/jair.953.

[25] Li Lisha, Jamieson Kevin, DeSalvo Giulia, Rostamizadeh Afshin, and Talwalkar Ameet. 2017. Hyperband: A novel bandit-based approach to hyperparameter optimization. J. Mach. Learn. Res. 18, 1 (2017), 6765–6816.

[26] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. Advances in neural information processing systems, 30.

[27] Ribeiro, Marco & Singh, Sameer & Guestrin, Carlos. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. 97-101. 10.18653/v1/N16-3020.

[28] Deepak Kumar Sharma, Jahanavi Mishra, Aeshit Singh, Raghav Govil, Gautam Srivastava, Jerry Chun-Wei Lin, (2022) *Explainable Artificial Intelligence for Cybersecurity, Computers and Electrical Engineering*, Volume103,108356,ISSN 00457906, https ://doi.org/10.1016/j.compeleceng.2022.108356.

# Low-Light Image Enhancement using Retinex-based Network with Attention Mechanism

Shaojin Ma[1], Weiguo Pan[2]*, Nuoya Li[3], Songjie Du[4], Hongzhe Liu[5], Bingxin Xu[6], Cheng Xu[7], Xuewei Li[8]*

Beijing Key Laboratory of Information Service Engineering, Beijing Union University, China[1, 2, 3, 4, 5, 6, 7, 8]
College of Robotics, Beijing Union University, Beijing Union University, China[1, 2, 3, 4, 5, 6, 7, 8]

*Abstract*—**Images in low-light conditions typically exhibit significant degradation such as low contrast, color shift, noise and artifacts, which diminish the accuracy of the recognition task in computer vision. To address these challenges, this paper proposes a low-light image enhancement method based on Retinex. Specifically, a decomposition network is designed to acquire high-quality light illumination and reflection maps, complemented by the incorporation of a comprehensive loss function. A denoising network was proposed to mitigate the noise in low-light images with the assistance of images' spatial information. Notably, the extended convolution layer has been employed to replace the maximum pooling layer and the Basic-Residual-Modules (BRM) module from the decomposition network has integrates into the denoising network. To address challenges related to shadow blocks and halo artifacts, an enhancement module was proposed to be integration into the jump connections of U-Net. This enhancement module leverages the Feature-Extraction- Module (FEM) attention module, a sophisticated mechanism that improves the network's capacity to learn meaningful features by integrating the image features in both channel dimensions and spatial attention mechanism to receive more detailed illumination information about the object and suppress other useless information. Based on the experiments conducted on public datasets LOL-V1 and LOL-V2, our method demonstrates noteworthy performance improvements. The enhanced results by our method achieve an average of 23.15, 0.88, 0.419 and 0.0040 on four evaluation metrics - PSNR, SSIM, NIQE and GMSD. Those results superior to the mainstream methods.**

*Keywords—Low-light image enhancement; decomposition network; FEM attention mechanism; denoising network; detail enhancement*

## I. INTRODUCTION

In the field of computer vision, low-light image enhancement has perennially been a focal point of research. Images acquired under low light conditions are often affected by problems like light weakening, increased noise and loss of detail, resulting in degraded image quality and blurred image content. The identified limitations exert a detrimental impact on the efficacy of computer vision applications, presenting challenges across various scenarios, including object detection [1], driverless driving [2], medical imaging. Moreover, these deficiencies introduce inconvenience in routine image capture and sharing within everyday life.

Traditional image enhancement methods often rely on manually adjusting parameters such as brightness and contrast [3], which may not adapt well to changes in different scenes.

However, due to the inability to effectively and accurately capture the features of the image, as well as its complex textures, these methods can lead to over-enhancement or the presence of shadow blocks and halo artifacts in the enhanced image. In contrast, methods based on deep learning improve adaptability and generalization by automatically learning image features, making them particularly suitable for complex environments. Specifically, deep learning methods based on Retinex theory, which separate an image's illuminance and reflectance through a decomposition network, allow for detailed adjustments through reflectance recovery and illuminance adjustment networks. These methods then merge the enhanced reflectance and illuminance images to improve brightness, contrast, and maintain natural colors, making them especially suitable for image enhancement in low-light environments. However, the decomposition network in Retinex can be affected by uneven lighting, potentially leading to loss of image detail, especially in dark and highlight areas of the image, thereby affecting the naturalness and realism of the final enhancement effect.

Although there are currently various enhancement methods for low-light images, mainly focusing on improving image contrast, it is important to note that low-light images often contain a significant amount of noise, which can greatly affect the quality and clarity of the image. Many of the current denoising techniques are applied in the pre-processing and post-processing stages of the image. Denoising in the pre-processing stage can cause the image to become blurred, while applying denoising in the post-processing stage can lead to the amplification of noise. Therefore, in the process of enhancing low-light images, how to appropriately balance the suppression of noise with the preservation of image details becomes a key challenge.

To address the aforementioned issues, this paper presents three main contributions, as follows:

*1) In* this paper, a decomposition network is proposed to obtain illumination and reflection maps through the decomposition of the RM and IM modules, as well as a comprehensive loss function is advanced to maintain the overall structure and consistency of the decomposed images.

*2) A* denoising network was proposed to remove noise in low-light images, with the assistance of images' spatial information, noise can be efficiently diminished, preserving map details, and consequently elevating the overall quality of enhanced images.

*3) To* effectively mitigate shadow blocks and halo artifacts in low-light images, this paper introduces an enhancement network featuring the FEM attention mechanism, which can significantly improve the restoration of image details and textures, yielding clear and natural image results.

The rest of this article is organized as follows: Section II discusses the related work. The proposed approach is detailed in Section III. Section IV provides quantitative and qualitative evaluation the method's performance. In Section V, ablation experiments were carried out on the FEM module and the de-noising module, respectively.

## II. RELATED WORK

Over the past few decades, various conventional methods have been proposed to address the challenges of low-light image enhancement. Noteworthy among these are traditional image enhancement methods, specifically histogram equalization [4] and methods based on Retinex theory [5]. Among these methods, histogram equalization method is one of the earliest and most extensively utilized methods. It aims to enhance image contrast and brightness by redistributing the gray level of image pixel values. This approach exists the disadvantages such as the limitations of global processing and the sensitivity to noise, which can lead to unnatural effects and information loss.

To enhance the visual alignment of images with human perception, Land et al. [6] proposed Retinex theory. At the core of the theory lies the concept that an object's color is determined not only by the intensity of the reflected light but also by its ability to reflect light waves. However, this method exhibits limitations when applied to images with complex lighting conditions and strong contrasts. To address these issues, researchers introduced the Multiscale Retinex (MSR) algorithm [7], incorporating multi-scale Gaussian filters to more accurately estimate the illumination components at various scales, and suppressed the halo effect through weighted summation. Furthermore, in pursuit of preserving the natural and authentic visual characteristics of images, Gao et al. [8] proposed an improved Retinex algorithm based on the traditional Retinex algorithm, which introduced color correction and multiscale processing technology to improve image details at different scales more precisely, thereby improving the visual effect and quality of images. However, most RetineX-based methods can cause severe color distortion and struggle to effectively enhance images with relatively high dynamic range.

In recent years, the remarkable adaptive capabilities of deep learning in low-light image enhancement have established it as an effective method, contributing to the improvement of image quality and finding widespread application in various computer vision tasks. Numerous scholars have extended their efforts to constructing learning-based models based on Retinex theory. For instance, RetinexNet [9] integrates Retinex theory with deep convolutional neural networks, enhancing image contrast through brightness maps estimation and adjustment, with subsequent post-processing using Block-Matching and 3D Filtering (BM3D) for denoising. Zhang et al. [10] designed an efficient network based on Retinex theory to enhance low-light images. Lim et al. [11] introduced the Deep-Stacked Laplacian Restorer (DSLR), capable of recovering global brightness and local detail from the original input, achieving notable success in contrast improvement and noise reduction. Moreover, several non-Retinex-based methods have been proposed. Li et al. [12] developed LightenNet, a convolutional neural network employing a stacked sparse denoising autoencoder structure to learn the nonlinear transformation function for adaptive brightness and contrast enhancement in low light images. However, this method faces challenges in effectively addressing noise in low-light images conditions. Ma et al. [13] proposed a low-light image enhancement method based on a fast, flexible and robust strategy. The method combines adaptive enhancement and parameter adjustment to efficiently enhance images while preserving detail and quality. EnlightenGAN [14] employs an unsupervised deep learning approach within a Generative Adversarial Network (GAN) framework to address low-light image enhancement challenges. Depth-Aware Decomposition and Restoration Network (DA-DRN) [15] introduces a self-sensing depth Retinex network, directly restores degraded reflectance and preserves the detail information in the decomposition stage by using the dependence between reflectance and illumination pattern. Although these methods can significantly improve the brightness and contrast of images, challenges persist in noise removal and image detail recovery. Some methods may result in overly enhanced image, leading to potential distortions.

The essence of the Retinex method lies in the estimation of luminance and reflectance maps. Traditional methods, with their limited decomposition ability, often result in over-enhancement or under-enhancement. In contrast, the learning-based approach demonstrates enhanced decomposition results and effectively improves contrast. It is noteworthy that many learning-based methods primarily utilize spatial information from low-light images to generate high-quality normal-light images, often neglecting the recovery of detailed information. Therefore, the enhancement module proposed in this paper leverages the FEM attention module, embedding it into U-Net's jump connection to augment the network's learning capacity for meaningful features. This augmentation is achieved by integrating image features and spatial attention mechanisms in the channel dimension. This design not only utilizes spatial information to strengthen contrast but also prioritizes the recovering of finer detail from the image.

In low-light conditions, image quality is often constrained by the optical signal attenuation, resulting in a significant degradation of the signal-to-noise ratio and a notable increase in noise. Some methods use the denoising model as preprocessing methods for low-light image enhancement; however, this preprocessing result in the loss of details in the low-light image. Chen et al. [16] introduced a model employing two parallel CNN branches: one for extracting brightness information and the other for extracting residual noise information. Low-light Photo Denoising via a Diffusion Model (LPDM) [17] is a denoising method for low-light images that utilizes a diffusion process. Initially, low-light images are enhanced to improve their brightness and contrast, followed by noise reduction through a diffusion process.

Although this method is effective in reducing noise, it introduces certain issues, such as the loss of some image details. To mitigate the loss of detail information during the denoising process, this paper introduces a denoising network to suppress noise in the reflection map. However, eliminating noise by restraining high-frequency signals in reflectivity map may lead to the loss of inherent details.



Fig. 1. The framework of the proposed method.

## III. PROPOSED METHOD

As illustrated in Fig. 1, the whole pipeline is based on Retinex. Initially, a comprehensive loss function is incorporated into the decomposition network to obtain light and reflection maps. In order to mitigate the loss of detailed information, a denoising network is introduced to suppress the high frequency information in the reflection image. Finally, the FEM attention module is proposed in the enhancement network to process the light image, aiming to better preserve object details, create smoother color transitions, and yield a clearer, more natural image.

### A. Decomposition Module

The loss function of decomposition network consists of perception loss function, illumination consistency loss function and global consistency loss function.

$$L = \lambda_1 \cdot L_{perceptual} + \lambda_2 \cdot L_{consistency} + \lambda_3 \cdot L_{global} \quad (1)$$

The values of $\lambda_1$, $\lambda_2$, $\lambda_3$ and are 0.3, 0.4, and 0.3.

The perception loss function is designed to preserve the perceived quality of the image. Traditional pixel level loss function often falls short in accurately capturing the perceived quality of the image. Hence, we choose a pre-trained convolutional neural network to extract image features and subsequently calculate the feature difference between the generated low-light image and the target image.

$$\mathrm{L}_{perceptual} = \frac{1}{N} \sum_{k=1}^{N} \| \phi(I)_k - \phi(L)_k \|_2^2 \quad (2)$$

Where, the input low-light image is $I$, the target light map is $L$, $\phi(I)$ represents the feature map extracted by the pre-trained convolutional neural network (such as VGG16), and $N$ is the number of feature maps, $\|.\|_2$ stands the $L2$ norm.

The illumination consistency loss function ensures the structural information's consistency between the generated low-light image and the target image.

$$L_{perceptual} = \frac{1}{N} \sum_{k=1}^{N} \| \mathrm{G}(I)_k - \mathrm{G}(L)_k \|_2^2 \quad (3)$$

The $\mathrm{G}()$ represents the gradient of the image.

The global consistency loss function elevates the overall consistency of the generated low-light image and the target image.

$$L_{global} = \frac{1}{N} \sum_{k=1}^{N} \| \mathrm{mean}(I)_k - \mathrm{mean}(L)_k \|_2^2 \quad (4)$$

The $mean()$ represents the mean of the image.

### B. BRM Module

The BRM module (Fig. 2) adopts the concept of a residual network as a reference and comprises 5 convolution layers. The convolution kernel size is {1, 3, 3, 1}, and the corresponding number of the convolution kernel is {64, 128, 128, 64}. For the activation function, SELU is assigned to correspond to the convolution kernel 3, and LeakyReLu to the convolution kernel 1. Finally, a 64×1×1 convolution layer added to the jump junction.



Fig. 2. The framework of BRM.

## C. Denoising Module

The heavy noise in low-light images obscures essential details, structure, and other valuable information, burying useful features beneath irrelevant ones. It is these severe degradations that make the training process challenging for network learning and the recovery of useful features (such as details, structure, and corrected color information).

Previous methods often eliminate noise by suppressing high-frequency signals in the reflection map. However, those signals frequently encompass critical image details and texture information. The inhibition of high-frequency signals can lead to detail loss or blurring, resulting in visually smooth or less sharp enhanced image. In this paper, a denoising network is posed to remove noise in low-light images by considering spatial information from images. U-Net [18] has demonstrated excellent results in numerous computer vision tasks, and it is frequently employed in low-light image enhancement networks. Nevertheless, U-Net's use of multiple max pooling layers has resulted in loss of feature information. In our network, the maximum pooling layer is replaced by the extended convolution layer, enabling a broader context information range by increasing the receptive field size of the convolution kernel without reducing the feature map resolution. The BRM module from the decomposition network is integrated into the denoising network, with each sub-module of U-Net being replaced with BRM. In the denoising network's encoder part, subsampling is achieved by adding an average pooling layer with a pool core size and step size of 2 to the BRM module at the end. In the decoder part, up-sampling is realized by incorporating a deconvolution layer with convolution kernel size and step size of 2 to the BRM module. The spatial information in images often contains rich details and structural information. The denoising network proposed in this paper leverages spatial information for denoising, enhancing its ability to preserve details and resulting in clearer and more natural images after denoising.

In the denoising network, a novel loss function is proposed in this paper. By integrating the Mean Square Error (MSE) loss term with the smoothing loss term, the image's noise suppression and smoothing effect can be optimized simultaneously. The MSE loss term aids in minimizing the pixel-level difference between the real image and the low-light image, thereby mitigating the impact of noise. The smoothing loss item promotes the smooth characteristics of the low-light image, enhancing the clarity of edges and details. The specific process is as follows:

$$L = \lambda Lsmooth + Lmse \tag{5}$$

The tradeoff between smoothness and the difference at the pixel level in the denoising loss function can be controlled by adjusting the weighting factor $\lambda$ for the smoothing loss term, which defaults to 0.2.

$$Lmse = ||Igi - {}^{\square}Igi||^2 \tag{6}$$

Where, $Igi$ represents the grayscale image relative to the low-light image, and ${}^{\square}Igi$ represents the grayscale image relative to the real image.

$$Lsmooth = ||\nabla \tilde{}\, Igi||^2 \tag{7}$$

Where, ${}^{\square}Igi$ represents the image processed by the denoising network, and $\nabla$ represents the gradient operator, which is used to calculate the gradient of the image. In this paper, the Prewitte operator is arranged to represent the value of the gradient operator, which can be represented by the following two matrices:

$$G_x = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}, G_y = \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} \tag{8}$$

## D. Enhancement Module

After denoising, residual shadow blocks and halo artifacts persist in low light images. In this paper, an enhancement module is devised to remove these artifacts while improving the quality of low-light images. Notably, the generation of shadow blocks and halo artifacts can be attributed to the U-Net's jump connection, wherein severely degraded features are directly conveyed to the up-sampled stage by linking up-sampled features with previous down-sampled features, leading to the retention of degraded features.

Inspired by SENet [19] in image recognition, the FEM attention module is incorporated into U-Net's jump connection to enhance noise removal and facilitate detail recovery. This inclusion is particularly effective in eliminating shadow blocks and halo artifacts. The FEM attention module operates by integrating image features in the channel dimension, assigning higher weights to valuable features (such as the correct color, detail, and texture features). This enabling the network to better learn these crucial features, while assigning lower or zero weights to less important features (such as noise, distorted colors, shadow blocks, and halo artifacts) or even giving no weight at all.

The loss function of the enhancement module is shown as follows:

$$L_{Re} = L_{Ren-con} + \lambda L_{Re-per} \tag{9}$$

Where $\lambda$ is the weight used to balance different loss terms, the default value is 0.1. $L_{Ren-con}$ the content loss value obtained by calculating the absolute value of the difference between the enhanced image and the real enhanced image at each position of pixel, and adding the absolute value of all differences. The presented loss metric quantifies the holistic disparity between the generated enhanced image and the authentic enhanced image. $L_{Re-per}$ is the perception loss, which is gained by computing the square difference in the perception space between the generated enhanced image and the actual enhanced image, summing across all pixel positions. To maintain a balanced consideration of the various layers within the feature map, normalization is conducted by dividing the dimensions of the feature map.

Content loss is defined as follows:

$$L_{Re-con} = \sum_{i}^{N} |{}^{\square} S_{low} - S_{en}| \tag{10}$$

Where $i$ is the index of the pixel, $\square S_{low}$ is the generated enhanced image, $S_{en}$ is the real enhanced image. $\|$ is the absolute value symbol. $\sum$ means adding the difference of each pixel, that is, adding the difference between the generated enhanced image and the real enhanced image at each pixel position.

Perceived loss means:

$$LRe - per = \frac{1}{(CjHjWj)} * \Sigma \parallel \varphi j(\square S_{low}) - \varphi j(S_{en}) \parallel^2 \qquad (11)$$

Where $Cj$ represents the number of channels in the feature map of the *J-th* layer and the number of channels in the feature map of different layers used in the perception loss. $Hj$ is the height of the feature map of the *J-th* layer, which represents the height of the feature map of the different layers used in perception loss. $Wj$ represents the width of the feature map of the *J-th* layer, representing the width of the feature map of the different layers used in perception loss. $\varphi j$ represents a function that maps the image to the *J-th* layer feature map, which is used to extract the representation of the features of image on the perceptual space. $\square S_{low}$ represents the enhanced image generated. $S_{en}$ is a true enhanced image. $\parallel \cdot \parallel^2$ represent the square of the Euclidean norm of two vectors used to calculate the square of the difference in perceptual space between the generated and real enhanced images.

### E. FEM Module

In recent years, a multitude of attention modules have been proposed to incorporate learnable weights in information processing, facilitating dynamic adjustments and the assignment of significance to various parts of the input data. This approach draws inspiration from human perceive and cognitive processes, enabling models to concentrate on information that significantly contributes to a given task or problem. For instance, Hu et al. [19] propose a Squeeze-and-Excitation (SE) block, which effectively performs feature recalibration by modeling the interrelationships between channels. Recognizing the importance of positional relationships between pixels, Non-Local Network (NLNet) [20] explored a nonlocal operation to capture the interactions between any two positions, irrespective of their spatial distance. Subsequently, the Cross Partial attention Volume Transformer (CPVT) module introduces an attention mechanism that crosses partial channels and divides the input feature map into several subgroups, each encompassing a subset of the channels. The attention mechanism is applied to each subgroup, focusing exclusively on channel relationships within the current subgroup and not considering channels in other subgroups. This approach achieves a balance between computational efficiency and performance.

These methods prove advantageous in addressing complex tasks like object detection and scene segmentation. A notable example is Axial-Deeplab [21], a deep learning model designed for image segmentation tasks. It employs an Axial attention mechanism for processing large-scale images and incorporates a segmented attention mechanism to enhance the capture of relationships between objects. However, these approaches may exhibit limited impact on low-level tasks, such as image enhancement. To address this issue, we propose an FEM (Fig. 3) attention module in this paper. The module effectively removes shadow blocks and halo artifacts by integrating image features and spatial attention mechanisms in the channel dimension. Furthermore, optimized jump connections are introduced, enabling adaptive exploration of contrast information within the image and facilitating the recovery of potentially fine details in low-brightness areas.



Fig. 3. The framework of FEM.

The FEM module comprises two Conv+ReLU layers, AdaptiveAvgPool2d, AdaptiveMaxPool2d, an up-sampling module and a Dual Attention Module (DAM). Specifically, for an input feature graph X with dimensions $H \times W \times C$, AdaptiveAvgPool2d and AdaptiveMaxPool2d are employed to extract representative information. The average of these two operations generates a global information feature map with dimensions $1 \times 1 \times C$. Then, the feature map with global information undergoes amplified through up-sampling, and the number of channels is compressed using $1\times1$ Conv to obtain a global feature map with dimensions $H \times W \times C1$. Following this, a DAM module is introduced to extract global features from spatial and channel dimensions. The DAM consists of two input branches: channel attention branch, and spatial attention branch. For an input feature graph X with dimensions of $H \times W \times C$, the channel attention branch employs global average pooling and global maximum pooling to generate the global average pooling feature map Cavg and the global maximum pooling feature map Cmax in spatial dimension, respectively. Their purpose is to emphasize the information regions, and these results are combined to produce the output Fc $(R1\times1\times C)$ of the attention branch of the channel. The spatial attention branch aims to generate a space-based attention map. Similar to the channel attention branch, it computes Savg $(RH\times W \times1)$ and Smax $(RH\times W \times1)$ through global average pooling and maximum pooling in the channel dimension respectively. The output spatial attention map Fs $(RH\times W \times1)$ is then obtained through a convolution layer. Finally, Fc and FS are combined to rescale and optimize the global feature map, resulting in the output. The input feature map (encoding local information) and the optimized global feature map (encoding global information) are combined using the

concatenate function and the Conv+ReLU function to generate an output feature map with dimensions H × W × C.

## IV. EXPERIMENTAL

### A. Parameter Setting

The experiment was conducted using PyTorch 1.8.0, with network training confined to a 256×256 patch on a single NVIDIA GTX 1080Ti GPU. The batch size was set to 2 and a total of $1\times10^5$ iterations were performed. Data enhancement involved random horizontal and vertical flips. Employing the Adam optimizer, the initial learning rate was set to $10^{-4}$, gradually reduced to $10^{-6}$ through a cosine annealing strategy.

### B. Data Set

The low light image dataset comprises a curated collection of specialized images tailored for the examination and evaluation of image processing algorithms in low light conditions. Typically, these datasets contain images captured in settings with insufficient illumination, offering researchers a challenging assortment for the development and assessment of algorithms aimed at enhancing the quality of low-light images. Derived real-world low-light scenes, these datasets encompass derived environments, both indoor and outdoor, capturing a range of shooting conditions and objects. These images within these datasets commonly exhibit characteristics such as low contrast, indistinct light and dark details, and elevated noise levels.

In the experiment of this paper, we utilized the LOL-v1 [9] and LOL-v2-real [22] datasets. The LOL dataset consists of 500 pairs of images, each with low and normal illumination, totaling 1,000 images. The images are captured with the same camera in varying lighting conditions, these images span a diverse array of scenes, both indoor and outdoor, and cover various shooting conditions and subjects. The dataset's diversity ensures comprehensive coverage of low-light scenes, which enhance the generalization and robustness of the evaluation algorithm. The LOL v2 dataset serves as a sequel to LOL v1, consists of two pairs of training and validation images, involving the actual shot and composite-generated images. Specifically, LOL-v2 is divided into two subsets: LOL-v2-REAL and LOL-v2-synthetic. The former includes 689 pairs of low-light/normal-light images for training and 100 pairs for testing, primarily adjusted by modifying camera parameters like exposure time and ISO. The latter is generated on an illumination distribution analysis of RAW format images.

### C. Evaluation Index

*1) Subjective evaluation:* The method presented in this paper is compared with several advanced low-light image enhancement methods, including KIND [10], KIND+ + [23], NE [24], SCI [13] and RetinexNet [9]. Experiments are conducted using publicly available source code provided by the authors of these methods.

The results depicted in Fig. 4, The RetinexNet method overly smoothens details and even causes color deviations, making the image look unnatural. Moreover, the results of RetinexNet still contain a lot of noise. Although it uses BM3D

to remove noise from the decomposed reflectance component, it cannot clearly remove the noise. The main reason is that BM3D is designed to remove Gaussian noise with a fixed noise level. However, the noise in the reflectance component is more diverse and complex than Gaussian noise. Although KIND and KIND + + introduced a recovery network to recovery color and remove noise from reflected images, the results were still inconsistent. For instance, in the first row of the Fig. 4, the KIND++ enhanced image still displays color deviation when compared to the reference image. In the second row, KIND treats a small light source as noise and removes it. In the sixth row, it is evident that KIND has unevenly enhanced the image. In the third image, the enhanced KIND ++ image still has color bias compared to the reference image. The SCI based method produces visually appealing results, it carries some undesirable artifacts (such as white walls). In contrast, upon observing the visualization results, particularly the outline of the teddy bear in the first row of Fig. 4 and the texture details of the book in the third row, our proposed method outperforms in terms of enhancement, reduced noise, and more accurate detail recovery. Conversely, other methods yield a fuzzy recovery effect due to noise interference in low light images, with the recovered images retaining significant noise. By comparing the results of different methods in the fifth line of images, our proposed method performs superior performance in restoring color saturation, presenting a more natural and vivid color enhancement effect. Moreover, in contrast to the restored color of the window in the second row of Fig. 4 and the floor in the last row, it can be concluded that our proposed method excels in maintaining color fidelity, suppressing noise, and removing artifacts.



Fig. 4. Subjective comparison on the LOL-V1 dataset.

Fig. 5 illustrates the impact of the experiment detailed in this paper, along with comparisons to other experiments on the LOL V2 dataset. Although RetinexNet can enhance the low-light areas in images, it exhibits severe color distortion and halo artifacts (a significant amount of halo artifacts can be observed around the contours of the mountains in the second row of images). In contrast, our method effectively brightens the low-light areas reasonably while maintaining the realistic visibility of the result. Besides enhancing brightness, our method also successfully reduces color distortion and halo artifacts. SCI introduces and even amplifies noise after enhancement, and therefore suffers from severe noise

distortion and color degradation when brightening dark areas. Conversely, our experimental results have superior color-recovery performance. KIND tends to exhibit some under-enhanced areas and apparent color distortions. Although the KIND ++ can remove noise, excessive sharpening (the surface of the mountain in the second row and the intersection of branches in the sixth picture) may introduce other problems such as loss of detail, blurring, and causing the image to unnatural. Compared to all these methods, our method can effectively enhance the brightness and display details of the image while suppressing noise, and it presents the most abundant and reasonable color information.



Fig. 5.    Subjective comparison on the LOL-V2dataset.

*2) Objective evaluation:* According to Table I, our method achieves superior results, which manifests that the enhanced images obtained by this pipeline are more visually satisfactory.

To objectively evaluate the enhancement results, we employed four classical indicators, and the results are presented in Table I. Among these metrics, PSNR [25] is utilized to measure the noise level and degree of distortion between the original image and the processed image. A higher PSNR value indicates a closer result to the reference image at the pixel-level. SSIM [26] is employed to evaluate structural similarity, considering perceptual properties such as image structure and content. A higher SSIM value indicates a greater similarity in structure to the reference image. The method proposed in this paper achieves excels in both PSNR and SSIM, highlighting its advantages in lighting restoration and structural restoration. KIND and NE also achieved high PSNR values, indicating their effectiveness in restoring global illumination. GMSD [27] assesses image quality by comparing gradient amplitude difference between the original and the processed images. NIQE [28] quantifies the effect of the image enhancement algorithm by analyzing the statistical

characteristics of the image and generating a continuous value score. A lower NIQE score indicates higher quality detail, brightness, and tone, indicative of a more natural appearance devoid of artifacts or pseudo-details. Consistently, our method achieves superior performance in both PSNR and SSIM. The enhanced images generated by our pipeline exhibit a more natural and vivid appearance, showcasing enhanced global and local contrast.

TABLE I.        COMPARISON OF OBJECTIVE EVALUATION INDICATORS OF DIFFERENT MODELS

|  | KIND | KIND++ | NE | SCI | RetinexNet | Our |
|---|---|---|---|---|---|---|
| PSNR | 21.38 | 19.21 | 22.61 | 20.80 | 18.4 | 23.15 |
| SSIM | 0.85 | 0.79 | 0.82 | 0.72 | 0.62 | 0.88 |
| NIQE | 0.610 | 0.556 | 0.526 | 0.470 | 0.831 | 0.419 |
| GMSD | 0.060 | 0.106 | 0.091 | 0.063 | 0.137 | 0.040 |

### D. Ablation Experiment

*1) Comparison of the effectiveness of FEM module:* To validate the efficacy of the FEM module, a model was trained with the FEM module replaced by an ordinary convolution. Fig. 6 illustrates the impact of FEM module on image enhancement. The results indicate that the model incorporating the FEM module effectively preserve object details and achieves smoother color transition in the image. Notably, in the third and fourth images, the enhanced results closely approximate the real image, aligning with the characteristics of the natural landscape. Analysis of the data presented in Table I reveals a notable improvement when utilizing models with FEM module compared to those without. Specifically, the PSNR increases by 3.74 (=20.22-16.48) and SSIM increases by 0.18 (=0.81-0.63). These finding lead to conclusion that models incorporating FEM modules exhibit superior performance in image enhancement.



Fig. 6.    Ablation experiments to verify the effectiveness of the FEM model.

*2) The effectiveness of the denoising network:* To validate the efficacy of the denoising module, a comparison was made between a method trained without the denoising network and the original method. As depicted in the Fig. 7, the

experimental results employing the denoising network effectively remove the noise from the image, and the results after the removal of noise contain finer details and more vivid colors. *It is* illustrate*d* in the Table *II* that the method incorporating the denoising network improves the PSNR ratio by 2.88 (22.77-19.29) and the SSIM ratio by 0.17 (0.84-0.67). These results demonstrate the effectiveness of the denoising network in this experiment.



Fig. 7.    Ablation experiments to verify the effectiveness of the FEM module.

TABLE II.        RESULTS OF THE OBJECTIVE EVALUATION OF THE ABLATION EXPERIMENTS

| Different situations | PSNR | SSIM | NIQE | GMSD |
|---|---|---|---|---|
| No FEM module | 16.48 | 0.63 | 0.51 | 0.083 |
| FEM module | 20.22 | 0.80 | 0.46 | 0.066 |
| No denoising network | 19.29 | 0.67 | 0.58 | 0.073 |
| denoising network | 22.17 | 0.84 | 0.43 | 0.042 |

## V.    CONCLUSIONS

This paper proposes a low light image enhancement method based on Retinex. We propose an efficient network for decomposing a low light image, incorporating an innovative loss function that integrates perceptual, light consistency and global consistency to obtain high quality light and reflection maps. The decomposition network comprises a reflection image extraction module (RM) and an illumination image extraction module (IM). Additionally, we integrate a denoising network and an enhancement module to further improve image quality. Our method not only enhances image color smoothness, reduce artifacts, but also effectively remove noise to restore image details under low-light conditions. By employing the FEM attention module instead of the convolution layer, our method successfully preserves object details. This results in a smoother color transition, yielding clearer and more natural images. Experimental results demonstrate that the proposed method achieves significant performance improvement across various low-light scenes. When compared to other existing methods, our method excels in image enhancement and detail recovery, showcasing superior noise removal and artifact suppression. In future

work, we aim to extend the application of this method to other computer vision tasks, substantiating its versatility and performance advantages in different domains through comparisons with other state-of-the-art methods.

## REFERENCES

[1]    Shrikhande, Saachi, Siddhesh Borse, and Shripad Bhatlawande. Face Recognition Based Attendance System. No. 10070. EasyChair, 2023.

[2]    Li, X., Chen, S., Hu, X., & Wang, J. (2021). Autonomous driving using deep learning: A survey. IEEE Transactions on Intelligent Transportation Systems, 23(10), 3789-3813.

[3]    Ma, S., Pan, W., Liu, H., Dai, S., Xu, B., Xu, C., ... & Guan, H. (2023). Image Dehazing Based on Improved Color Channel Transfer and Multiexposure Fusion. Advances in Multimedia, 2023.

[4]    Pizer S M, Amburn E P, Austin J D, et al. Adaptive histogram equalization and its variations[J]. Computer vision, graphics, and image processing, 1987, 39(3): 355-368.

[5]    Abdullah-Al-Wadud, M., Kabir, M. H., Dewan, M. A. A., & Chae, O. (2007). A dynamic histogram equalization for image contrast enhancement. IEEE transactions on consumer electronics, 53(2), 593-600.

[6]    Land, Edwin H. "The retinex theory of color vision." Scientific american 237.6 (1977): 108-129.

[7]    Jobson, D. J., Rahman, Z. U., & Woodell, G. A. (1997). Multiscale Retinex for color image enhancement. IEEE Transactions on Image Processing, 6(7), 965-976.

[8]    Gao, Yuhang, Chuhao Su, and Zhaoheng Xu. "Color image enhancement algorithm based on improved Retinex algorithm." 2022 Asia Conference on Algorithms, Computing and Machine Learning (CACML). IEEE, 2022.

[9]    Wei, C., Wang, W., Yang, W., & Liu, J. (2018). Deep retinex decomposition for low-light enhancement. arXiv preprint arXiv:1808.04560.

[10]   Zhang, Yonghua, Jiawan Zhang, and Xiaojie Guo. "Kindling the darkness: A practical low-light image enhancer." Proceedings of the 27th ACM international conference on multimedia. 2019.

[11]   Lim, Seokjae, and Wonjun Kim. "DSLR: Deep stacked Laplacian restorer for low-light image enhancement." IEEE Transactions on Multimedia 23 (2020): 4272-4284.

[12]   Li C, Guo J, Porikli F, et al. LightenNet: A convolutional neural network for weakly illuminated image enhancement[J]. Pattern recognition letters, 2018, 104: 15-22.

[13]   Ma, Long, et al. "Toward fast, flexible, and robust low-light image enhancement." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

[14]   Jiang, Y., Gong, X., Liu, D., Cheng, Y., Fang, C., Shen, X., ... & Wang, Z. (2021). Enlightengan: Deep light enhancement without paired supervision. IEEE transactions on image processing, 30, 2340-2349.

[15]   X. Wei, X. Zhang, S. Wang, C. Cheng, Y . Huang, K. Yang, and Y . Li,"Da-drn: Degradation-aware deep retinex network for low-light image enhancement," arXiv preprint arXiv:2110.01809, 2021.

[16] Chen, C., Chen, Q., Xu, J., & Koltun, V. (2018). Learning to see in the dark. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 3291-3300).

[17] Panagiotou, Savvas, and Anna S. Bosman. "Denoising Diffusion Post-Processing for Low-Light Image Enhancement." arXiv preprint arXiv:2303.09627 (2023).

[18] Ronneberger, Olaf, Philipp Fischer, and Thomas Brox. "U-net: Convolutional networks for biomedical image segmentation." Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18. Springer International Publishing, 2015.

[19] Hu, Jie, Li Shen, and Gang Sun. "Squeeze-and-excitation networks." Proceedings of the IEEE conference on computer vision and pattern recognition. 2018.

[20] Wang, X., Girshick, R., Gupta, A., & He, K. (2018). Non-local neural networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 7794-7803).

[21] Wang, Huiyu, et al. "Axial-deeplab: Stand-alone axial-attention for panoptic segmentation." European conference on computer vision. Cham: Springer International Publishing, 2020.

[22] Yang, W., Wang, W., Huang, H., Wang, S., & Liu, J. (2021). Sparse gradient regularized deep retinex network for robust low-light image enhancement. IEEE Transactions on Image Processing, 30, 2072-2086.

[23] Zhang, Y., Guo, X., Ma, J., Liu, W., & Zhang, J. (2021). Beyond brightening low-light images. International Journal of Computer Vision, 129, 1013-1037.

[24] Jin, Yeying, Wenhan Yang, and Robby T. Tan. "Unsupervised night image enhancement: When layer decomposition meets light-effects suppression." European Conference on Computer Vision. Cham: Springer Nature Switzerland, 2022.

[25] Huynh-Thu, Quan, and Mohammed Ghanbari. "Scope of validity of PSNR in image/video quality assessment." Electronics letters 44.13 (2008): 800-801.

[26] Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing, 13(4), 600-612.

[27] Xue, W., Zhang, L., Mou, X., & Bovik, A. C. (2013). Gradient magnitude similarity deviation: A highly efficient perceptual image quality index. IEEE transactions on image processing, 23(2), 684-695.

[28] Mittal, Anish, Anush Krishna Moorthy, and Alan Conrad Bovik. "No-reference image quality assessment in the spatial domain." IEEE Transactions on image processing 21.12 (2012): 4695-4708.

# Double Branch Lightweight Finger Vein Recognition based on Diffusion Model

Zhiyong Tao[1], Yajing Gao[2], Sen Lin[3]

School of Electronic and Information Engineering, Liaoning Technical University, Huludao, China 125105[1, 2]
School of Automation and Electrical Engineering, Shenyang Ligong University, Shenyang, China 110159[3]

*Abstract*—**Aiming at the problems of high complexity, insufficient global information extraction and easy overfitting in finger vein recognition, a finger vein recognition method based on diffusion model is proposed. Firstly, finger vein images are generated according to the dataset by diffusion model, which is used to prevent overfitting; secondly, a streamlined convolutional neural network is used to form a two-branch lightweight backbone network with an improved multi-head self-attention mechanism, which can effectively reduce the complexity of the model; and finally, in order to maximally extract the image's overall information, the convolution is used to merge the extracted local and global features, and the recognition results are output. The algorithm can reach a maximum recognition rate of 99.78% on multiple datasets, while the number of references is only 2.15M, which further reduces the complexity of the algorithm while maintaining a high accuracy compared to other novel finger vein recognition algorithms as well as lightweight convolutional neural network models. As the first attempt in this field, it will provide new ideas for future research work.**

*Keywords—Finger vein recognition; convolution neural network; diffusion model; multi-head self-attention mechanism; lightweight network*

## I. INTRODUCTION

Lately, the focus of researchers on finger vein recognition technology has intensified, attributed to its exceptional security and precision. Since the vein network of each individual is hidden under the skin, finger vein-based biometrics has a massive advantage in live identification. As a developing technology, finger vein recognition-based biometrics is far from flawless. Various internal and external factors can impact the performance of finger vein verification. These factors include: Lighting Conditions, Finger Placement Angle, and Uniform illumination. Therefore, high-precision and high-robustness algorithms are essential for feature extraction, recognition, or verification of finger vein images. A typical finger vein recognition process includes image acquisition, preprocessing, feature extraction, and matching. In the finger vein image acquisition process, a finger vein image acquisition device consisting of an image sensor and an infrared light source is used. Preprocessing of the acquired finger vein images is carried out to facilitate the subsequent feature extraction process. Suppressing noise, improving image contrast, and performing data augmentation are common image preprocessing methods.

Feature extraction in finger vein recognition involves two main categories: traditional recognition methods and deep learning methods. Traditional recognition methods for feature extraction can be further classified into three distinct categories: template-based methods [1], representation-based methods [2], and feature-based [3][4][5] learning methods. These methods require manual labeling of parameters, depend on image quality, and have cumbersome recognition steps. Compared with machine learning methods, deep learning [6][7][8][9] based methods can achieve more stable recognition results by acquiring more profound image features through Convolutional Neural Networks (CNN). Therefore, some researchers proposed deep learning-based finger vein recognition methods. For example, Radzi et al. [10] proposed a CNN-based finger vein recognition method, Fang et al. [11] proposed a lightweight two-channel network to improve the verification of finger veins by extracting the mini-region of interest (ROI), Zhang et al. [12] proposed Domain Adaptation Finger Vein Network (DAFVN) improve the final recognition result by extracting illumination invariant features in the image and reducing the effect of light on the recognition result. Recently, researchers proposed the Vision Transformer (ViT) [13] method, which has attracted widespread attention in deep learning. Compared with CNN, ViT focuses more on global features and has shown excellent performance in several domains. In addition, researchers have proposed some improved methods, such as Liu et al. [14] proposed Swin Transformer, which obtains global and local features by constructing hierarchical feature maps and sliding windows, with better experimental results but high model complexity, and Peng et al. [15] proposed Parallel Network Architecture, which makes use of convolution and Multi-Head Self-Attention (MHS) mechanism. Head Self-Attention (MHSA) to extract local and global features in parallel, which improves the network performance but is ineffective for small datasets. Based on the advantages of Transformer, researchers started applying it to finger vein recognition. Huang et al. [16] proposed Finger Vein Transformer (FVT) model for recognition, which achieves multi-scale feature extraction by reducing the number of tokens layer by layer, but exploiting the Transformer increases the complexity and computation at the same time.

To enhance the efficiency of recognizing finger veins, some researchers have introduced data enhancement techniques to make the model better adapt to finger vein images in various scenarios. Yang et al. [17] proposed Finger Vein Representation Using the Generative Adversarial Networks (FV-GAN) model, which was the first time GAN was in the field of finger vein recognition. Choi et al. [18] proposed a Conditional Generative Adversarial Network (CGAN) to recover blurred images. They used a deep convolutional neural

network for the finger vein images—a convolutional neural network for finger vein image recognition. Hou et al. [19] proposed a ternary classifier, GAN, for generating training data to improve the learning ability of the CNN classifier. Although high-quality images can be generated using GAN networks, they require more extensive databases and may be unstable during training.

Diffusion Model (DM) [20] is a recently emerged deep generative model for high-quality image generation, which is rapidly evolving and is widely used in tasks such as text-to-image generation, image-to-image generation, and video image generation. Also, one of the data enhancement methods, the diffusion model has a more straightforward training process and generates higher-quality images compared to GAN networks. Jonathan et al. [20] proposed the Denoising Diffusion Probabilistic Model (DDPM), which is used for image generation tasks, generating the image quality is higher than other generation models such as GAN. Robin [21] et al. proposed Latent Diffusion Models (LDM), which achieve image generation by introducing a cross-attention conditioning mechanism, which significantly improves the training and sampling efficiency without degrading its quality.

Summarizing the above research methods, the existing algorithms for finger vein recognition generally have high complexity, recognition accuracy needs to be improved, and the training process is unstable, so a two-branch lightweight finger vein recognition model (FV-DM) based on the diffusion model is designed to achieve finger vein image generation

using the diffusion model to solve the problem of overfitting due to the small finger vein dataset. The CNN and the improved E-MSHA module are used to extract image features in parallel with the dual-branching in the feature extraction process to avoid the problem of low accuracy caused by insufficient feature extraction, while the diffusion model is used in the finger vein recognition process in order to explore a new way of finger vein recognition. The comprehensive experiments on the self-constructed dataset and three public datasets show that FV-DM all achieve better recognition results, as well as lower model parameters and computational complexity, shorter recognition time, and lower Equal Error Rate (EER).

The remainder of this paper is organized as follows. Section 2 introduces our modelling approach and explains how it works. Section 3 describes the experiments conducted to validate the performance of the model. In Section 4, we summarize the paper and make suggestions for future work.

## II. METHODOLOGY

### A. Diffusion Model

Recently, the diffusion model as a generative model has received more and more attention from researchers due to its powerful image generation ability. As shown in Fig. 1, the diffusion model is mainly divided into the forward noise addition process and the reverse denoising process. The solid line indicates the forward noise addition process and the dashed line indicates the reverse denoising process.



Fig. 1. Network structure diagram of diffusion model.

*1) Forward noise addition process:* The forward noise addition process uses Gaussian noise to gradually add noise to the input image, generating a series of noise samples $x_0, x_1,...,x_T$ until the image becomes a pure noise image. Assuming that $q(x_0)$ is the probability distribution of the real image, $q(x_t | x_{t-1})$ represents the probability distribution of the current image $x_t$ obtained by adding noise to the previous step image $x_{t-1}$ in the forward noise addition process, and the mathematical expressions for each step of the process of adding Gaussian noise are shown in (1):

$$q(x_t | x_{t-1}) = N(x_t | \sqrt{1-\beta_t}\, x_{t-1}, \beta_t I) \quad (1)$$

Where $\beta_t$ is the diffusivity and $t$ varies with time. The formula is expressed as a mean $\mu_t = \sqrt{1-\beta_t}\, x_{t-1}$ with a variance Gaussian $\sigma_t^2 = \beta_t$ distribution. If the final image $x_T$ is obtained through $x_0$, the whole process can be regarded as a Markov chain from $t=1$ to the moment $t=T$, as shown in Eq. (2):

$$q(x_{0:T}) = q(x_0)\prod_{t=1}^{T} q(x_t | x_{t-1}) \quad (2)$$

In the forward noise addition process, Eq. (1) can be expressed as by the simplified way in literature [23]:

$$x_t = \sqrt{1-\beta_t}\, x_{t-1} + \sqrt{\beta_t}\, z_{t-1} \qquad (3)$$

Where $z_{t-1}$ denotes the noise at moment $t-1$. The $x_t$ at any moment is obtained from the original image $x_0$ with the formula shown in the following equation:

$$\alpha_t = 1 - \beta_t \qquad (4)$$

$$\overline{\alpha_t} = \prod_{i=1}^{t} \alpha_i \qquad (5)$$

$$\begin{aligned} x_t &= \sqrt{\alpha_t}\, x_{t-1} + \sqrt{1-\alpha_t}\, z_{t-1} \\ &= \sqrt{\overline{\alpha_t}}\, x_0 + \sqrt{1-\overline{\alpha_t}}\, z_t \end{aligned} \qquad (6)$$

$$q(x_t \mid x_0) = N\,(x_t | \sqrt{\overline{\alpha_t}}\, x_0, (1-\overline{\alpha_t})I) \qquad (7)$$

$z_t$ is denoted as a Gaussian distribution satisfying $N\,(0, I)$

*2) Reverse denoising process:* The reverse denoising process is also known as the inverse diffusion process. The main purpose is to gradually predict the target image $x_0$ from the purely noisy image $x_T$, i.e., to derive the $x_{t-1}$ distribution from $x_t$, which can be transformed into what is shown in Eq. (8) by using Bayes' formula:

$$q(x_{t-1} \mid x_t) = q(x_t \mid x_{t-1})\frac{q(x_{t-1})}{q(x_t)} \qquad (8)$$

According to the forward noise addition process, $q(x_t \mid x_{t-1})$ is known, and for $q(x_{t-1})$ and $q(x_t)$, it can be

solved by adding the known condition $x_0$, as shown in the following equation:

$$q(x_{t-1} \mid x_t, x_0) = q(x_t \mid x_{t-1}, x_0)\frac{q(x_{t-1} \mid x_0)}{q(x_t \mid x_0)} \qquad (9)$$

In the process of reverse denoising, the features in the input noisy image are predicted by the neural network, and this paper chooses U-Net as the model for noise prediction. U-Net is a U-shaped network structure, which consists of downsampling on the left side, upsampling on the right side, and cross-layer connections. The downsampling reduces the size of the feature map through the convolution operation and reduces the computational cost. The upsampling gradually restores the feature map to its original size through the inverse convolution operation, and the cross-layer connection is used to splice the features between the downsampling and the upsampling, which can effectively integrate the features of different levels of the image. For normalisation, Group Normalization (GN) is chosen. Finally, for the downsampling and upsampling operations in U-Net, the convolution with a step size of 2 and the inverse convolution are chosen, respectively. The specific structure is shown in Fig. 2.

### B. Design of the Network Model

Inspired by DDPM [20], a finger vein recognition network based on a diffusion model is designed. It is shown in Fig. 3



Fig. 2. U-Net Structure diagram.

Fig. 3. Structure diagram of diffusion model.

The network structure contains two main parts: the image generation part and the image feature extraction part. Firstly, the diffusion model is used to achieve image generation by forward noise addition process and reverse denoising process. Then, the generated image is passed into the feature extraction network along with the actual image, and the Residual [22] module and the E-MHSA module extract the global and local features of the image, respectively, and stitch the features after extraction, which allows for better fusion of the features and further improves their expressive ability. The fused features are passed through the convolution module to achieve the extraction of deeper features. As shown in the figure, the convolution module consists of an ordinary convolution of size, an ordinary convolution of size, and a maximum pooling layer and is stacked twice in the feature extraction process to extract image features more comprehensively.

### C. Residual Structure

The model uses a Residual module and an improved E-MHSA module for local and global feature extraction in the early stage of feature extraction. The Residual module consists of an inverted residual structure, which can effectively reduce the computational cost while extracting the local features of the image. The specific structure is shown in Fig. 4.



Fig. 4. Residual structure diagram.

The inverted residual structure contains two ordinary convolutions, a DW convolution, a Dropout layer and a jump connection. The input information is first increased by $1\times1$ size ordinary convolution, then the image size is transformed by DW convolution with a convolution kernel size of $3\times3$, and finally the number of channels is decreased by $1\times1$ size ordinary convolution, and the Dropout is used to randomly discard the features to prevent the parameter from relying too much on the training data and the phenomenon of overfitting. Finally, the output of the Dropout layer is added with the result of the jump join to complete the output of the information. The jump connection is mathematically defined as:

$$H(x) = F(x) + x \qquad (10)$$

Where $F(\cdot)$ is a function containing convolution, pooling, and modified linear unit operations, $x$ inputs the feature map, and $H(x)$ is the output of the inverted residual structure. The inclusion of jump connections in the inverted residual accelerates the convergence of the network and improves the generalization of the model.

### D. E-MHSA Structure

Since MHSA has a strong ability to capture low-frequency signals, which are used to provide global information, the enhanced E-MHSA module is used in this paper for global feature extraction. Compared to the traditional MHSA, E-MHSA incorporates average pooling operation and down-sampling before the computation of the attention mechanism in order to reduce the computational cost and achieve a more efficient and lightweight deployment. As shown in Fig. 5, the E-MHSA module is similar to the Transformer Block in ViT, which first captures the low-frequency signals through E-MHSA with the following formula:

$$\text{E-MHSA}(x) = concat(SA(x_1), SA(x_2), ..., SA(x_h))W^0 \qquad (11)$$

Where $x = [x_1, x_2, ..., x_h]$ denotes the division of input feature $x$ into multiple heads in the channel dimension and $h$ is the number of heads divided. In this paper, we take 8 as the number of heads for the attention mechanism. $SA(\cdot)$ is the computational formula for the attention mechanism, and the formula is as follows:

$$SA(x) = Attention(X \cdot W^Q, P_s(X \cdot W^K), P_s(X \cdot W^V)) \quad (12)$$

where $P_s$ represents the average pooling operation with step size $s$.



Fig. 5. Structure diagram of E-MHSA.

### III. EXPERIMENTS AND ANALYSES

#### A. Presentation of Datasets

The experiments were conducted on three public datasets, FV-USM [24], SDUMLA-HMT [25], THU-FVFDT2 [26], and with a self-constructed dataset, FV-SIPL, which were divided in a 2:1 ratio, except for the THU-FVFDT2 dataset, in which the training and test sets were equally divided. The data information is shown in Table I.

TABLE I. DATA INFORMATION FROM FOUR DATASETS

| Dataset | Total number of categories | Total image count | Total training sets | Total test sets |
|---|---|---|---|---|
| FV-USM | 492 | 5904 | 3936 | 1968 |
| SDUMLA-HMT | 636 | 3816 | 2544 | 1272 |
| THU-FVFDT2 | 610 | 1220 | 610 | 610 |
| FV-SIPL | 108 | 1296 | 864 | 432 |

*1) FV-USM:* The dataset was provided by Universiti Teknologi Malaysia and contained finger vein images from 123 volunteers, with 12 images captured from each of the four fingers of each volunteer. Therefore, the whole dataset covers a total of 492 finger categories and 5904 images. The size of each of these images is 640×480pixels.

*2) SDUMLA-HMT:* The dataset was provided by Shandong University, which contains finger vein images of 106 volunteers, and 6 images were collected for each index, middle, and ring finger of each volunteer's hands the whole dataset covers a total of 636 finger categories and 3816 images, where each image size is 320×240pixels.

*3) THU-FVFDT2:* The dataset was provided by Tsinghua University and contained finger vein images of 610 volunteers. Finger vein images were collected twice for each volunteer, with a total of 1220 images, each with a size of 200×100pixels.

*4) FV-SIPL:* This dataset was made by the Signal and Information Processing Laboratory of Liaoning University of Engineering and Technology by using infrared finger vein acquisition sensors to collect finger vein images from 27 volunteers. Among them, 12 images were acquired for each of the four fingers of each volunteer, and the whole dataset covered 108 finger categories and 1296 images in total. The size of each image is 176×415 pixels.

#### B. Image Preprocessing

In order to facilitate the subsequent process of image feature extraction, preprocessing operations are performed on the image. Taking the FV-USM dataset as an example, the main processes are shown in Fig. 6(a) to (d) below.



(a) Raw image     (b) ROI Extraction

(c) Normalised image    (d) Diffusion modelling to generate images    (e) GAN-generated images

Fig. 6. Image preprocessing process.

For the original images in the dataset, first of all, through the ROI extraction operation, to reduce the interference of irrelevant information on the recognition results, and then carry out image normalisation, pass the normalised images into the diffusion model, and set the number of iterations in the training process to be 10000, and the time $T$ to be 1000. the same parameter settings are carried out on the commonly used generative model GAN, and it can be seen through Fig. 6(d) and Fig. 6(e) that the images generated by the diffusion model are clearer and show more similar image features to the original image, so the use of diffusion model is chosen as the data enhancement method in FV-DM.

#### C. Experimental Environment and Parameter Settings

The experiments were conducted under the Linux operating system using PyTorch1.7 framework, and the graphics card used for training and testing was GeForce RTX 3090. The learning rate was set to 0.001, the batch size was set to 16, and Stochastic Gradient Descent (SGD) was chosen as the optimiser, where the momentum was set to 0.9. The input size of finger veins was uniformly adjusted to 224×224pixels, and the final experimental results were obtained by training iterations 100 times.

#### D. Evaluation Indicators

In order to evaluate the performance and advantages of the model, metrics such as Accuracy, Equal Error Rate, Average Processing Time for a Single Image, Number of Parameters, and Floating Point Operations (FLOPs) are selected for evaluation. Accuracy rate, as one of the commonly used metrics in finger vein recognition, can reflect the ability of the model to correctly identify different categories of samples in the entire dataset. The formula for accuracy rate is shown in (13):

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (13)$$

Where *TP* denotes the number of correct positive sample predictions, *TN* denotes the number of correct negative sample predictions, *FN* denotes the number of incorrect negative sample predictions, and *FP* denotes the number of incorrect positive sample predictions. In image recognition tasks, the EER value is usually used as an indicator to evaluate the good or bad performance of the model, which is determined by the False Acceptance Rate (FAR) and the False Rejection Rate (FRR). The formulas for FAR and FRR are shown below:

$$FAR = \frac{FP}{FP + TN} \qquad (14)$$

$$FRR = \frac{FN}{TP + FN} \qquad (15)$$

Wherein the number of samples for incorrect acceptance and incorrect rejection is defined by a predetermined threshold. When the threshold of matching is greater than the preset threshold, it is determined to be incorrectly accepted, and vice versa is determined to be incorrectly rejected. The value when *FRR* and *FAR* are equal is the equal error rate. The equal error rate reflects the overall performance of the recognition method. The smaller the value of equal error rate, the better the performance of the recognition method.

*E. Comparison Experiment*

In order to verify the effectiveness of the FV-DM method, it is compared with the classical Transformer network models: the VIT-B, Swin-T, Conformer-B, Next-ViT and the lightweight CNN network model EfficientNetV2. The recognition accuracy results of the different methods on the datasets are shown in Table II. The results in the table show that the methods proposed in this paper achieve the best recognition results on all four datasets. Bolding indicates the best results and underlining indicates the second best results. In addition to the accuracy comparison, the average processing time, number of parameters and FLOPs of individual images for the different methods were also

compared, as shown in Table III. In terms of the average processing time for a single image, MobileNetV2 is 2.27ms, which is 0.53ms faster than FV-DM, which is due to the MSHA contained in FV-DM. Other than that, FV-DM outperforms the other methods.

TABLE II. RECOGNITION ACCURACY OF DIFFERENT METHODS ON FOUR DATA SETS (UNIT: %)

| Method | FV-USM | SDUMLA-HMT | THU-FVFDT2 | FV-SIPL |
|---|---|---|---|---|
| VIT-B[13] | 58.67 | 63.66 | 59.55 | 76.28 |
| Swin-T[14] | 95.0 | 93.33 | 76.01 | 95.12 |
| Conformer-B[15] | 99.0 | 98.67 | 96.64 | 99.33 |
| Next-ViT[27] | 98.56 | 99.0 | 98.87 | 99.53 |
| EfficientNetV2[28] | 98.10 | 98.07 | 97.78 | 98.20 |
| MobileNetV2[22] | 98.12 | 99.0 | 98.32 | 99.0 |
| ResNet101[29] | 98.33 | 98.34 | 98.21 | 99.0 |
| FV-DM(Our) | **99.67** | **99.66** | **99.10** | **99.78** |

TABLE III. COMPARISON OF EVALUATION INDEX RESULTS OF DIFFERENT METHODS

| Method | Time/ms | Parameters/M | FLOPs/G |
|---|---|---|---|
| VIT-B | 11.30 | 103.03 | 16.88 |
| Swin-T | 7.21 | 28.27 | 4.37 |
| Conformer-B | 7.15 | 96.63 | 21.01 |
| Next-ViT | 3.52 | 31.76 | 5.79 |
| EfficientNetV2 | 3.49 | 21.46 | 2.90 |
| MobileNetV2 | **2.27** | 3.50 | 0.33 |
| ResNet101 | 7.61 | 44.55 | 7.84 |
| FV-DM(Our) | 2.80 | **2.15** | **0.19** |

In this paper, four datasets are used to compare different recognition methods, including VIT-B, Swin-T and Conformer-B. The results are shown in Fig. 7.



Fig. 7. Comparing the equal error rates of different methods.

On the SDUMLA-HMT dataset, the equal error rate of FV-DM is slightly higher than that of MobileNetV2, but except for that, FV-DM maintains the lowest equal error rate, which indicates that the FV-DM method has excellent performance in finger vein recognition, and it can be used as an effective recognition method. Compared with other methods, the FV-DM method has higher accuracy and better robustness, so it has a wide range of application prospects in practical applications.

TABLE IV.       RECOGNITION ACCURACY OF DIFFERENT METHODS ON PUBLIC DATASETS (UNIT: %)

| Method | FV-USM | SDUMLA-HMT | THU-FVFDT2 |
|---|---|---|---|
| Merge CNN[30] | 96.15 | 89.99 | — |
| DS-CNN[31] | — | 98.00 | 89.00 |
| Semi-PFVN[32] | 94.67 | 96.61 | — |
| LFVRN_CE[33] | 98.58 | 97.75 | — |
| DGLFV[34] | — | 99.25 | — |
| CMrFD[35] | 98.33 | 98.92 | — |
| FVT | **99.73** | 97.90 | 90.66 |
| TFHFT-DPFNN[36] | — | 98.00 | — |
| CNNs[37] | 97.95 | — | — |
| Coding SchemeA[38] | 99.59 | 95.91 | — |
| FV-GAN | — | — | 98.52 |
| Triplet-classifier GAN | 99.66 | 99.53 | — |
| FV-DM(Our) | 99.67 | **99.66** | **99.10** |

FV-DM is compared with novel finger vein models in recent years, and the results are shown in Table IV. Among them, FV-DM obtained the highest recognition accuracy on both public datasets, SDUMLA-HMT and THU-FVFDT2. The recognition accuracy on the FV-USM dataset is lower than that of the FVT method by 0.06%, but it is higher than that of FVT on the SDUMLA-HMT and THU-FVFDT2 datasets by 1.77% and 9.03%, respectively. Therefore, from the overall results, FV-DM recognition results are better.



(a) Recognition accuracy curves for the four datasets.



(b) Loss curves for the four datasets.

Fig. 8.   Recognition accuracy and loss curve of FV-DM on four datasets.

By comparing the novel finger vein recognition algorithms in recent years, FV-DM has better performance in terms of recognition accuracy, recognition time, complexity, etc. the recognition accuracy versus test loss curves of FV-DM on the four datasets are shown in Fig. 8.

*F. Ablation Experiment*

Ablation experiments were conducted in order to better validate the effectiveness of the modules in each part of the network. Under the premise that the rest of the conditions remain unchanged, modules such as Residual, E-MHSA, convolutional module and diffusion model are added to the network sequentially, and the accuracy rate is tested with the FV-SIPL dataset as an example, and the experimental results are shown in Table V. From the table, it can be seen that the accuracy rate increases step by step after the modules are added. Residual and E-MHSA need to be further fused after extracting the local and global features, respectively, to achieve a more comprehensive feature extraction, so the accuracy rate is increased by 42.92% after adding the convolution module compared with the previous one. The introduction of the diffusion model can achieve intra-class enhancement of the data and avoid the overfitting problem, so the accuracy is further improved after adding the diffusion model, which verifies the correctness of the conjecture.

TABLE V.       ACCURACY COMPARISON ON THE FV-SIPL DATASET (UNIT: %)

| Residual | E-MHSA | Convolution module | Diffusion model | Accuracy |
|---|---|---|---|---|
| √ | | | | 41.40 |
| √ | √ | | | 55.58 |
| √ | √ | √ | | 98.50 |
| √ | √ | √ | √ | **99.78** |

IV.     CONCLUSION

Aiming at the finger vein recognition process that does not fully consider the global features of the image, is easy to overfit, and has other problems, this paper proposes a two-

branch lightweight finger vein recognition method based on the diffusion model. Firstly, the diffusion model is used to generate finger vein images to expand the finger vein dataset. Secondly, a two-branch network composed of a convolutional neural network and improved E-MHSA is used to extract global and local features from the expanded dataset. Then, the extracted global and local features are fused by the convolutional module, and the image features are further extracted. Finally, the recognition results are output, and the effectiveness of the method is verified on multiple datasets at the same time. Experiments show that the method in this paper can improve the recognition performance while keeping the computational cost small. In future work, the application of the diffusion model in finger vein recognition will be explored deeply to seek more possibilities.

ACKNOWLEDGMENT

REFERENCES

[1] Yang L, Yang G P, Yin Y L, et al. Finger Vein Recognition With Anatomy Structure Analysis[J]. IEEE Transactions on Circuits and Systems for Video Technology, 2018, 28: 1892-1905.

[2] Lu Y, Yoon S, Wu S Q, et al. Pyramid Histogram of Double Competitive Pattern for Finger Vein Recognition[J]. IEEE Access, 2018, 6: 56445-56456.

[3] Kang W X, Lu Y T, Li D J, et al. From Noise to Feature: Exploiting Intensity Distribution as a Novel Soft Biometric Trait for Finger Vein Recognition[J]. IEEE Transactions on Information Forensics and Security, 2019, 14: 858-869.

[4] Yang L, Yang G P, Xi X M, et al. Finger Vein Code: From Indexing to Matching[J]. IEEE Transactions on Information Forensics and Security, 2019, 14: 1210-1223.

[5] Liu H Y, Yang G P, Yang L, et al. Anchor-Based Manifold Binary Pattern for Finger Vein Recognition[J]. Science China Information Sciences, 2019, 62: 1-16.

[6] Chen Y, Dai X, Chen D, et al. Mobile-former: Bridging mobilenet and transformer[C]//Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022: 5270-5279.

[7] Zhang Q, Yang Y B. Rest: An efficient transformer for visual recognition[J]. Advances in neural information processing systems, 2021, 34: 15475-15485.

[8] Lee S H, Lee S, Song B C. Vision transformer for small-size datasets[J]. arXiv preprint arXiv:2112.13492, 2021.

[9] Tan M, Le Q V. Mixconv: Mixed depthwise convolutional kernels[J]. arXiv preprint arXiv:1907.09595, 2019.

[10] Radzi S A, Khalih-hani M, Bakhteri R. Finger-Vein Biometric Identification Using Convolutional Neural Network[J]. Journal of Signal Processing, 2016, 24:1863-1878.

[11] Fang Y X, Wu Q X, Kang W X. A Novel Finger Vein Verification System Based on Two-Stream Convolutional Network Learning[J]. Neurocomputing, 2018, 29: 100-107.

[12] Zhang Z J, Zhong F, Kang W X. Study on Reflection-Based Imaging Finger Vein Recognition[J]. IEEE Transactions on Information Forensics and Security, 2021, 17: 2298-2310.

[13] Dosovitskiy A, Beyer L, Kolesnikov A, et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale[J]. arXiv preprint arXiv:2010.11929, 2020.

[14] Liu Z, Lin Y T, Cao Y, et al. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows[C]//Proceedings of the IEEE/CVF international conference on computer vision. 2021: 10012-10022.

[15] Peng Z L, Huang W, Gu S Z, et al. Conformer: Local Features Coupling Global Representations for Visual Recognition [C]//Proceedings of the IEEE/CVF international conference on computer vision. 2021: 367-376.

[16] Huan J D, Luo W J, Yang W L, et al. FVT: Finger vein transformer for authentication[J]. IEEE Transactions on Instrumentation and Measurement, 2022, 71: 1-13.

[17] Yang W M, Hui C Q, Chen Z Q, et al. FV-GAN: Finger vein representation using generative adversarial networks[J]. IEEE Transactions on Information Forensics and Security, 2019, 14: 2512-2524.

[18] Choi J, Noh K J, Cho S W, et al. Modified conditional generative adversarial network-based optical blur restoration for finger-vein recognition[J]. IEEE Access, 2020, 8: 16281-16301.

[19] Hou B, Yan R. Triplet-classifier GAN for finger-vein verification[J]. IEEE Transactions on Instrumentation and Measurement, 2022, 71: 1-12.

[20] Ho J , Jain A , Abbeel P .Denoising diffusion probabilistic models[J]. arXiv preprint arXiv:2006.11239, 2020.

[21] Rombach R, Blattmann A., Lorenz D, et al. High-resolution image synthesis with latent diffusion models [C]//Conference on Computer Vision and Pattern Recognition (CVPR), 2021: 10674-10685.

[22] Luo C. Understanding Diffusion Models: A Unified Perspective[J]. arXiv preprint arXiv: 2208.11970, 2022.

[23] Sandler M, Howard A, Zhu M, et al. Mobilenetv2: Inverted residuals and linear bottlenecks [C]//Proceedings of the IEEE conference on computer vision and pattern recognition. New York, 2018: 4510-4520.

[24] Asaari M S M, Suandi S A, Rosdi B A. Fusion of band limited phase only correlation and width centroid contour distance for finger based biometrics[J]. Expert Systems with Applications, 2014, 41(7): 3367-3382.

[25] Yin Y, Liu L, Sun X. SDUMLA-HMT: A multimodal biometric database [C]//Biometric Recognition: 6th Chinese Conference. Beijing, 2011: 260-268.

[26] Yang W, Qin C, Liao Q. A database with ROI extraction for studying fusion of finger vein and finger dorsal texture [C]//Biometric Recognition: 9th Chinese Conference. Shenyang, 2014: 266-270.

[27] Li J, Xia X, Li W, et al. Next-vit: Next generation vision transformer for efficient deployment in realistic industrial scenarios[J]. arXiv preprint arXiv:2207.05501, 2022.

[28] Tan M, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks [C]//International conference on machine learning. New York, 2019: 6105-6114.

[29] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition [C]//Proceedings of the IEEE conference on computer vision and pattern recognition. New York, 2016: 770-778.

[30] Zhao D, Ma H, Yang Z, et al. Finger vein recognition based on lightweight CNN combining center loss and dynamic regularization[J]. Infrared Physics & Technology, 2020, 105: 103221.

[31] Shaheed K, Mao A, Qureshi I, et al. DS-CNN: A pre-trained Xception model based on depth-wise separable convolutional neural network for finger vein recognition[J]. Expert Systems with Applications, 2022, 191: 116288.

[32] Chai T, Li J, Prasad S, et al. Shape-driven lightweight CNN for finger-vein biometrics[J]. Journal of Information Security and Applications, 2022, 67: 103211.

[33] Zhong Y, Li J, Chai T, et al. Different Dimension Issues in Deep Feature Space for Finger-Vein Recognition [C]//Chinese Conference on Biometric Recognition. Cham, 2021: 295-303.

[34] Tao Z, Wang H, Hu Y, et al. DGLFV: Deep Generalized Label Algorithm for Finger-Vein Recognition[J]. IEEE Access, 2021, PP(99):1-1.

[35] Shen J, Liu N, Xu C, et al. Finger vein recognition algorithm based on lightweight deep convolutional neural network[J]. IEEE Transactions on Instrumentation and Measurement, 2021, 71: 1-13.

[36] Muthusamy D, Rakkimuthu P. Trilateral Filterative Hermitian feature transformed deep perceptive fuzzy neural network for finger vein verification[J]. Expert Syst. Appl. 2022, 196: 116678.

[37] Zhao D , Ma H , Yang Z , et al. Finger vein recognition based on lightweight CNN combining center loss and dynamic regularization[J]. Infrared Physics & Technology, 2020, 105(8):103221.

[38] Ren H, Sun L, Guo J, et al. Finger vein recognition system with template protection based on convolutional neural network[J]. Knowl. Based Syst. 2021, 227, 107159.

# An Ensemble Approach to Question Classification: Integrating Electra Transformer, GloVe, and LSTM

Sanad Aburass[1], Osama Dorgham[2], Maha Abu Rumman[3]

Department of Computer Science, Maharishi International University, Fairfield, Iowa, USA[1, 3]
Prince Abdullah bin Ghazi Faculty of Information and Communication Technology, Al-Balqa Applied University Al-Salt, Jordan[2]
School of Information Technology, Skyline University College, University City of Sharjah Sharjah, United Arab Emirates[2]

*Abstract*—**Natural Language Processing (NLP) has emerged as a critical technology for understanding and generating human language, with applications including machine translation, sentiment analysis, and, most importantly, question classification. As a subfield of NLP, question classification focuses on determining the type of information being sought, which is an important step for downstream applications such as question answering systems. This study introduces an innovative ensemble approach to question classification that combines the strengths of the Electra, GloVe, and LSTM models. After being tried thoroughly on the well-known TREC dataset, the model shows that combining these different technologies can produce better outcomes. For understanding complex language, Electra uses transformers; GloVe uses global vector representations for word-level meaning; and LSTM models long-term relationships through sequence learning. Our ensemble model is a strong and effective way to solve the hard problem of question classification by mixing these parts in a smart way. The ensemble method works because it got an 80% accuracy score on the test dataset when it was compared to well-known models like BERT, RoBERTa, and DistilBERT.**

*Keywords—Ensemble learning; long short term memory; transformer models; Electra; GloVe; TREC dataset*

## I. INTRODUCTION

There are many areas where machine learning has completely changed how we solve problems. These include healthcare, banking, and natural language processing [1], [2], [3]. It has made it possible for computers to learn from data on their own, making choices, predicting trends, and even finding patterns that are too complicated for humans to understand. NLP is the study of how computers and people use language. With the rise of machine learning, big steps forward have been made in NLP, especially in areas like mood analysis, machine translation, and summary [4], [5], [6], [7]. One of the most important things that natural language processing does is sort questions into groups. In the real world, this job is very important for many things, such as search engines, virtual helpers like Siri or Google Assistant, and customer service bots. Question sorting that is done right can lead to more accurate and useful answers, which improves the service these apps can provide. Think about a medical robot that can correctly classify a health question and give a possibly life-saving answer, or a virtual tourist helper that can tell the difference between questions about food and questions about historical sites. It's not just handy that the good effects happen; they often have big effects [8], [9], [10]. However, the complexity of human language, which includes subtleties in syntax, meaning, and pragmatics, makes it very hard to get very accurate question classification [11], [12]. Support Vector Machines, Random Forests, and other machine learning models have been used for this, but new developments in deep learning and transformer models like BERT, RoBERTa, and ELECTRA have shown that they work even better than expected [13]. These models are very good at understanding the meanings and contexts of words and sentences, which is a key part of question classification [1], [14], [15], [16] and [17]. Here, we show a new method that combines three strong tools: the ELECTRA model for contextual embeddings based on transformers; Global Vectors for Word Representation (GloVe) for creating semantically rich word vectors; and Long Short-Term Memory (LSTM) networks for capturing sequence dependencies. The Text REtrieval Conference (TREC) dataset, which is a common standard for question classification tasks, is used to train and test our ensemble model. The main thing that our work adds is that we combine several different but useful techniques in a way that makes them work better together than current best models at classifying questions.

This study is organized into the following taxonomy: Section II starts by doing a full literature review of earlier work that looked at question categorization and related ensemble methods, Section III shows a full explanation of the method used is given, which includes the ELECTRA model, GloVe embeddings, and LSTM networks, Section IV presents the proposed approach, Section V describes how the experiment was set up, what the results were, and why we came to the conclusions we did, and in Section VI, we talk about the results, the limits, and the opportunities for more study.

## II. LITERATURE REVIEW

### A. Previous Work

In NLP, question categorization has been a major area of study for twenty years, with many researchers working on it. Over the years, techniques in this area have changed a lot, from simple machine learning methods to the most advanced deep learning models used today. Support Vector Machines (SVM) and other well-known machine learning methods were used in the early stages of this study. For example, Zhang and Lee used SVMs to sort questions [18].

Deep learning methods came out as machine learning got better. These made models more stable. Kalchbrenner et al. were the first to use convolutional neural networks (CNNs) to tag words with questions and put them into groups. After that, scientists studied Recurrent Neural Networks and various types of them, such as Long Short-Term Memory networks. After Zhou et al. used LSTMs well to find the long-term connections in question replies, they came up with some hopeful results [19].

When language models like BERT, RoBERTa, and ELECTRA came out, they were the next big step forward in the field of NLP. A lot of natural language processing jobs, like question classification, were done better by these transformer-based systems. Devlin et al. created BERT and showed that it could record context-rich embeddings [18]. While Liu et al. worked on RoBERTa and Clark et al. worked on ELECTRA, they pushed the limits of efficiency [20], [21].

Individual models have worked well on their own, but group methods have become popular as a way to combine the different strengths of these models. Vaswani et al. suggested a group that combined transformers and LSTMs, which showed a big improvement in performance compared to using just one model [22]. However, ensemble methods that are specifically made for question classification have not been widely used. This points to an interesting area for future study.

The role of word embeddings, especially GloVe, is another part of this changing environment. When Pennington et al. first presented GloVe, it quickly became a mainstay in many NLP tasks, such as question classification [23].

Before they come up with a new type of feature based on question patterns, Nguyen and Le look at lexical, syntactical, and semantic features. The writers came up with a way to choose features that would work for different types of questions. They used the TREC dataset and Support Vector Machines (SVM) for classification to show that their plan worked [24].

Chotirat and Meesad use two datasets—TREC-6 (English) and a Thai speech dataset—to test different machine learning models. The combined CNN-BiLSTM model did better than the other models, according to the findings. These results show that deep learning methods, especially mixed models, can improve the accuracy of question sorting in a lot of languages. The addition of Part-of-Speech tagging was a key factor in this speed boost [25].

The real-world data that Madabushi et al. give show that their system works better. When fine-grained question classification is paired with deep learning models, they show big improvements in how well the answers are chosen. The new taxonomy and object recognition system worked better than earlier models, showing that their way works. These results show how important it is to include question classification in deep learning systems for jobs like answer choice [26].

### B. Rationale for the Proposed Approach

Combining Electra, GloVe, and LSTM in a new way, we describe a new ensemble method for question classification,

this method was chosen because it can work well with others to help with the complex nature of understanding questions, with its transformer-based structure, Electra is great at handling complex language tasks and fully understanding their context, GloVe adds to this by providing detailed word-level meaning models that describe the complexity of how language is used, and LSTM helps by correctly simulating long-term relationships in text, which is very important for understanding how questions are asked in a certain order. These models work together to get around the problems that separate models like BERT and RoBERTa have, especially when it comes to handling complicated question forms and changing contexts. As you can see from our positive test results, our approach uses the strengths of each model to make question sorting more accurate and faster. This combination not only makes performance measures better, but it also makes it possible to analyze questions in a more detailed and full way, which is a big step forward in natural language processing.

Different modeling strategies have their own pros and cons, and there hasn't been much research on how to combine them into a single model for question classification, our work introduces a new ensemble method that combines ELECTRA, GloVe, and LSTM, the objective is to create a new style for grouping questions into different categories.

### III. BACKGROUND

This section provides a comprehensive overview of the primary components of our ensemble model: the ELECTRA model, GloVe word embeddings, and LSTM networks.

### A. ELECTRA

ELECTRA (Efficiently Learning an Encoder that Classifies Token Replacements Accurately) is a transformer-based model developed for natural language processing tasks, proposed by researchers at Google Research in 2020, ELECTRA uses a novel approach to training known as Replaced Token Detection [17].

Traditional transformer models, such as BERT [1], utilize masked language modeling as a pre-training task, where some percentage of the input tokens are masked and the model is trained to predict the original tokens. ELECTRA, on the other hand, introduces a different mechanism. It consists of two parts: a generator and a discriminator. The generator is a small masked language model that suggests replacements for some of the tokens in the input. The discriminator is then tasked with predicting whether each token in the sequence was replaced by the generator or not.

This training mechanism can be described with the following steps:

*1)* The generator G, a small BERT-like model, is used to replace some tokens in the input sequence.

*2)* The discriminator D, a larger BERT-like model, then attempts to predict for each position whether it contains the original token or a replacement.

The main advantage of this approach is that it allows for the entire input sequence to be utilized during pre-training, as

opposed to just a small masked portion, making the training process more efficient and effective.

*B. GloVe*

GloVe is an unsupervised learning algorithm developed by the Stanford NLP Group for obtaining vector representations for words. The primary idea behind GloVe is that the co-occurrence statistics of words in a corpus capture a significant amount of semantic information [23]. To construct the GloVe representations, the following steps are carried out:

*1)* A global word-word co-occurrence matrix is constructed from the corpus, where each element `$X_{ij}$` represents the frequency with which word `i` appears in the context of word `j`.

*2)* The objective of GloVe is then to learn word vectors such that their dot product equals the logarithm of the words' probability of co-occurrence.

Mathematically, this is represented as:

$$V_i \cdot V_j = \log(P(i|j)) \qquad (1)$$

where $V_i$ and $V_j$ are the word vectors for words i and j, and $P(i|j)$ is the probability of i appearing in the context of j.

*C. LSTM*

LSTM networks are a type of recurrent neural network (RNN) architecture [27], specifically designed to address the vanishing gradient problem of traditional RNNs and to better capture dependencies in sequential data [28]. In an LSTM, the hidden state $h_t$ is updated via a series of gating mechanisms:

*1)* The input gate $i_t$ determines how much of the new input will be stored in the cell state.

*2)* The forget gate $f_t$ decides the extent to which the previous cell state $c_{(t-1)}$ is maintained.

*3)* The output gate $o_t$ controls how much of the internal state is exposed to the external network.

The state update equations are as follows:

$$i_t = \sigma(W_{ii}.x_t + b_{ii} + W_{hi}.h_{(t-1)} + b_{hi}) \qquad (2)$$

$$f_t = \sigma(W_{if}.x_t + b_{if} + W_{hf}.h_{(t-1)} + b_{hf}) \qquad (3)$$

$$g_t = \tanh(W_{ig}.x_t + b_{ig} + W_{hg}.h_{(t-1)} + b_{hg}) \qquad (4)$$

$$o_t = \sigma(W_{io}.x_t + b_{io} + W_{ho}.h_{(t-1)} + b_{ho}) \qquad (5)$$

$$c_t = f_t * c_{(t-1)} + i_t * g_t. \qquad (6)$$

$$h_t = o_t * \tanh(c_t). \qquad (7)$$

Here, $\sigma$ represents the sigmoid function, tanh is the hyperbolic tangent function, * denotes element-wise multiplication, and `.` represents matrix multiplication. The variables W and b are the learnable weights and biases, respectively, of the LSTM.

By employing these gating mechanisms, LSTMs can effectively learn what information to keep or forget over long sequences, making them particularly efficient for tasks involving sequential data.

The combination of ELECTRA, GloVe, and LSTM in our ensemble model aims to leverage the efficient pre-training and high performance of ELECTRA, the rich semantic information encapsulated by GloVe embeddings, and the sequence modeling capabilities of LSTM. This synergistic integration seeks to enhance the performance of question classification tasks by capturing the semantics, context, and sequence information embedded in the questions [29], [30], [31].

IV. PROPOSED APPROACH

The proposed approach is designed to amalgamate the capabilities of multiple state-of-the-art language models and embeddings, namely Electra, GloVe, and LSTM, to enhance the classification performance on questions from the TREC dataset. The architecture employs a dual-branch neural network with each branch responsible for processing a different type of embedding—Electra for one and GloVe for the other. Subsequent to this, LSTM layers are applied to the concatenated embeddings, leading to the final classification output.

*A. Source of Data*

Based on the TREC question classification dataset, which has text-based questions and their related broad terms like "location," "person," etc., the experiment was carried out.

*B. Text Standardization*

The TensorFlow method tf.strings.lower() was used to change all of the raw text strings to lowercase.

*C. Tokenization and Sequence Padding*

There were two different tokenization processes for the raw texts: one was made for Electra and the other was made for GloVe. Through padding, a set sequence length of 512 was kept.

*D. Architectural Elements: In-Depth Exploration Electra Sub-model: Capturing Contextual Relationships*

Electra is the main tool used to find complex and detailed trends in searches. When it comes to Electra, the discriminator is very good at figuring out what a sign means in relation to its surroundings. This is very important for question classification because questions often have clues in the environment that help with classification. For instance, the use of "when" or "what year" could mean a question about time, which Electra is very good at spotting.

*E. GloVe Sub-model: Leveraging Global Statistical Information*

GloVe is useful because it can gather global statistical features of words based on data about how often they appear together. GloVe, unlike local context, records long-term ties like synonyms or similar ideas, which can be very helpful for finding the right questions. Electra can understand how the words in a question work together in complex ways, but GloVe takes it a step further by understanding the bigger language features of the words used.

*F. LSTM Layers: Accounting for Sequential Dependencies*

After integration, LSTM networks are used to find the sequence-based relationships in the incoming text. Questions naturally go in a certain order, with "wh" words like "who," "what," and "where" at the beginning and a subject or object at the end. Figuring out this process can often help you figure out what the question is really asking. These gates in LSTMs help them successfully capture long-term relationships, which makes them perfect for this job. The two LSTM layers, which have 256 and 128 units, are set up to add another level of abstraction and pick up more complex models.

*G. Classification Layer: Mapping to Categories*

The last Dense layer is a classifier that turns the complicated feature representations learned by the layers above into classification choices that can be used. In this case, a softmax activation function is used because the job is classified. There are 6 units in this layer, and each one represents a different type of question in the TREC dataset. The softmax function makes sure that the result can be understood as odds that add up to 1. It's easy to put each question into one of the six broad groups this way.

*H. Model Synergy: The Bigger Picture*

It is important to note that the architecture is not just a random group of techniques; it is a carefully put together set of techniques that are meant to work around the weaknesses and make the most of the strengths of each part. Electra gathers background, GloVe adds breadth, and LSTMs record how things change over time. These steps work together to make a complete plan for learning how to classify questions.

To put it simply, each design part was carefully chosen and put together in a way that makes a whole model that can change, understand, and do a great job of question classification.

## V. EXPERIMENTAL RESULTS

*A. Experimental Setup*

To thoroughly test how well our suggested ensemble model, Ensemble Electra + GloVe+LSTM, worked, we set up our tests on Google Colab Pro and used its GPU features to make the computations go faster. We put our ensemble model up against Electra and other cutting-edge language models [17], BERT [32], RoBERTa [33], and DistilBERT [34].

*B. Mathematical Overview of Models*

*1) ELECTRA:* Electra employs a discriminative training mechanism, where the model learns to distinguish between "real" and "fake" tokens in a sentence. Formally, for a given input $X = [x_1, x_2,…, x_n]$, a generator $G$ proposes replacements $x_i$ for masked tokens, and a discriminator $D$ estimates the probability $P(D(x_i) = 1| X)$ that each token is real. The objective is to minimize $-\log(D(x_i))$ for real tokens and $-\log(1 -D(\tilde{x}_i))$ for fake tokens.

*2) BERT:* BERT uses a masked language model (MLM) for pre-training, where a certain percentage of input tokens are masked. The model aims to predict these masked tokens based on their context. Mathematically, for an input sequence X, the loss L is calculated as $-\log P(x_i | X_{-i}; \theta)$, where $\theta$ are the model parameters.

*3) RoBERTa:* RoBERTa extends BERT but employs dynamic masking and removes the next-sentence prediction objective. Its objective function remains similar to BERT, focusing on masked token prediction.

*4) DistilBERT:* DistilBERT is a distilled version of BERT, trained to approximate BERT's output. For each token $x_i$ in the input X, the model aims to minimize the difference between its output $O(x_i)$ and that of BERT $B(x_i)$, typically using the Kullback-Leibler divergence.

*C. Evaluation Metrics*

We used several metrics to evaluate the performance of each model: Loss, Accuracy, Precision, Recall, and F1 Score.

*1) Loss:* Represents the error between predicted and actual labels. Lower values are better.

*2) Accuracy:* Measures the ratio of correctly predicted samples to the total samples.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \tag{8}$$

*3) Precision:* Indicates the percentage of positive identifications that were actually correct.

$$Precision = \frac{TP}{TP+FP} \tag{9}$$

*4) Recall:* Shows the percentage of actual positives that were identified correctly.

$$Recall = \frac{TP}{TP+FN} \tag{10}$$

*5) F1 Score:* Harmonic mean of precision and recall, a balance between the two.

$$F1\ Score = 2 * \frac{Precision*Recall}{Precision*+Recall} \tag{11}$$

Where: TP: True Positive, TN: True Negative, FP: False Positive and FN: False Negative.

*D. Results*

Our ensemble model, which is a combination of Electra, GloVe, and LSTM, outperformed all other models. The superior performance of our ensemble approach can be attributed to the complementary strengths of the constituent models. Electra, with its discriminator-generator setup, excels at understanding the context of the language. GloVe, on the other hand, captures semantic relationships between words by considering the global word-word co-occurrence statistics. LSTM effectively handles the sequence nature of the language data. Together, they give a complete approach to text classification and lead to great results on the TREC question classification task. This experimental evidence supports our theory that an ensemble of models can significantly improve question classification task performance over standalone models. By leveraging the strengths of each model, we were able to achieve superior results, showing that our proposed ensemble approach works. The results of the experiments are shown in Tables I and II and Fig. 1, 2, 3, 4 and 5.

TABLE I.        THE ACCURACY AND MSE OF THE MODELS

| Model | Train Accuracy | Test Accuracy | Train MSE | Test MSE |
|---|---|---|---|---|
| **Ensemble Electra + GloVe+LSTM** | **0.999** | **0.8** | **0.001** | **1.51** |
| Electra [17] | 0.229 | 0.188 | 5.055 | 5.44 |
| BERT [32] | 0.224 | 0.13 | 3.628 | 4.128 |
| RoBERTa [33] | 0.254 | 0.16 | 3.608 | 4.108 |
| Distilbert [34] | 0.239 | 0.145 | 3.628 | 4.128 |

TABLE II.        THE PRECISION, RECALL AND F1 SCORE OF THE MODELS

| Model | Train Precision | Test Precision | Train Recall | Test Recall | Train F1 Score | Test F1 Score |
|---|---|---|---|---|---|---|
| **Ensemble Electra+ GloVe+LSTM** | **0.999** | **0.8** | **0.999** | **0.8** | **0.999** | **0.8** |
| Electra | 0.052 | 0.035 | 0.229 | 0.188 | 0.085 | 0.0595 |
| BERT | 0.05 | 0.016 | 0.224 | 0.13 | 0.082 | 0.029 |
| RoBERTa | 0.08 | 0.046 | 0.254 | 0.16 | 0.112 | 0.059 |
| Distilbert | 0.065 | 0.031 | 0.239 | 0.145 | 0.097 | 0.044 |



Fig. 1.   Models accuracies.

Fig. 2.    Models precision.



Fig. 3.    Models recall.



Fig. 4.    Models F1 score.



Fig. 5.    Models mean squared error,

## VI. RESULTS AND DISCUSSION

All of the comparison data show that the Ensemble Electra + GloVe + LSTM model does better than all of the evaluation factors. This victory isn't just a small step forward; it's a huge step forward from solo ideas.

### A. Generalization and Overfitting

The ensemble model's ability to transfer from training data to test data is one of the most interesting results. With a training accuracy of 0.999 and a test accuracy of 0.8, the ensemble model shows that it can successfully apply learned patterns to data that it has never seen before. This even result shows that the model does not overfit, which is a common problem in machine learning [35].

### B. Error Analysis

The ensemble model stays ahead when it comes to Mean Squared Error (MSE). The model's predictions were very close to the real results, with a training MSE of 0.001 and a test MSE of 1.51. Standalone models, like Electra, BERT, and others, have much higher MSE values on both the training and test sets, which means they make more mistakes when making predictions.

### C. Precision, Recall, and F1 Score

The ensemble model also keeps its high scores in the F1 score, precision, and recall. A high accuracy score means that the ensemble model correctly finds relevant examples on a big scale, and a high recall score means that the model catches most of the relevant events. The F1 score, which is a fair way to measure precision and recall, shows that the model is well-balanced.

### D. Comparative Model Analysis

Although RoBERTa seems to do better than the other models that work by themselves, it is still not as good as the ensemble model. The ensemble model is the only one that can get Electra's understanding of context, GloVe's semantic depth, and LSTM's sequential reading all at the same time.

### E. Synergistic Strength

The enormous success of the ensemble model shows that combining parts that are similar to other cutting-edge models can create something new. For the TREC question answering test, it does very well because it knows data very well in both its specific and broad parts. The ensemble model does a great job of categorizing questions, and these results suggest that it could also help with other natural language processing issues.

## VII. CONCLUSION

In conclusion, our results show that an ensemble model with Electra, GloVe, and LSTM does a better job of classifying questions than other models on the TREC dataset. We tested our ensemble method against other advanced models like BERT, RoBERTa, and DistilBERT and found that it regularly did better than them. It achieved high accuracy, precision, recall, F1 score, and lower mean squared error. Electra, GloVe, and LSTM all have properties that work well together in the ensemble model. Combining different models and methods into ensemble methods, which we found, can lead to big performance gains, making them a reliable and

effective way to handle difficult tasks like question categorization. Even though these results are positive, we know that there is still room for improvement and adjustment. For instance, different groupings of ensembles and model designs could be looked into, along with more advanced training methods. In the future, researchers may look into how this ensemble method can be used to solve other natural language processing problems besides question classification. Overall, this study adds to the progress being made in natural language processing and lays the groundwork for more research and development of group methods in question categorization and other areas.

## VIII. CONFLICT OF INTEREST

The authors declare that there is no conflict of interest in this paper.

Declaration of Generative AI and AI-assisted technologies in the writing process

During the preparation of this work the authors used Quillbot in order to proofread the manuscript. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

## REFERENCES

[1] R. K. Kaliyar, 'A Multi-layer Bidirectional Transformer Encoder for Pre-trained Word Embedding: A Survey of BERT', in 2020 10th International Conference on Cloud Computing, Data Science & Engineering (Confluence), IEEE, Jan. 2020, pp. 336–340. doi: 10.1109/Confluence47617.2020.9058044.

[2] S. M. Rezaeinia, R. Rahmani, A. Ghodsi, and H. Veisi, 'Sentiment analysis based on improved pre-trained word embeddings', Expert Syst Appl, vol. 117, pp. 139–147, Mar. 2019, doi: 10.1016/j.eswa.2018.08.044.

[3] S. Aburass, A. Huneiti, and M. B. Al-Zoubi, 'Classification of Transformed and Geometrically Distorted Images using Convolutional Neural Network', Journal of Computer Science, vol. 18, no. 8, pp. 757–769, 2022, doi: 10.3844/jcssp.2022.757.769.

[4] T. Zhang, A. M. Schoene, S. Ji, and S. Ananiadou, 'Natural language processing applied to mental illness detection: a narrative review', NPJ Digit Med, vol. 5, no. 1, p. 46, Apr. 2022, doi: 10.1038/s41746-022-00589-7.

[5] K. S. Kalyan, A. Rajasekharan, and S. Sangeetha, 'AMMUS : A Survey of Transformer-based Pretrained Models in Natural Language Processing', Aug. 2021, [Online]. Available: http://arxiv.org/abs/2108.05542.

[6] S. Aburass, O. Dorgham, and M. A. Rumman, 'Comparative Analysis of LSTM and Ensemble LSTM Approaches for Gene Mutation Classification in Cancer', in 2023 IEEE International Conference on Machine Learning and Applied Network Technologies (ICMLANT), IEEE, Dec. 2023, pp. 1–6. doi: 10.1109/ICMLANT59547.2023.10372993.

[7] M. Fisher, O. Dorgham, and S. D. Laycock, 'Fast reconstructed radiographs from octree-compressed volumetric data', Int J Comput Assist Radiol Surg, vol. 8, no. 2, pp. 313–322, Mar. 2013, doi: 10.1007/s11548-012-0783-5.

[8] S. Aburass and O. Dorgham, 'Performance Evaluation of Swin Vision Transformer Model using Gradient Accumulation Optimization Technique', Jul. 2023, [Online]. Available: http://arxiv.org/abs/2308.00197.

[9] J. Al Shaqsi, O. Drogham, and S. Aburass, 'Advanced machine learning based exploration for predicting pandemic fatality: Oman dataset', Inform Med Unlocked, vol. 43, p. 101393, 2023, doi: 10.1016/j.imu.2023.101393.

[10] S. AbuRass, A. Huneiti, and M. B. Al-Zoubi, 'Enhancing Convolutional Neural Network using Hu's Moments', International Journal of Advanced Computer Science and Applications, vol. 11, no. 12, pp. 130–137, Dec. 2020, doi: 10.14569/IJACSA.2020.0111216.

[11] F. A. Acheampong, H. Nunoo-Mensah, and W. Chen, 'Transformer Models for Text-based Emotion Detection: A Review of BERT-based Approaches Network Optimisations View project Quantitative Medical Imaging View project Transformer Models for Text-based Emotion Detection: A Review of BERT-based Approaches'. [Online]. Available: https://www.researchgate.net/publication/348740926.

[12] O. Dorgham, I. Al-Mherat, J. Al-Shaer, S. Bani-Ahmad, and S. Laycock, 'Smart System for Prediction of Accurate Surface Electromyography Signals Using an Artificial Neural Network', Future Internet, vol. 11, no. 1, p. 25, Jan. 2019, doi: 10.3390/fi11010025.

[13] S. Aburass, O. Dorgham, and J. Al Shaqsi, 'A Hybrid Machine Learning Model for Classifying Gene Mutations in Cancer using LSTM, BiLSTM, CNN, GRU, and GloVe', Jul. 2023, [Online]. Available: http://arxiv.org/abs/2307.14361.

[14] S. Aburass, A. Huneiti, and M. B. Al-Zoubi, 'Classification of Transformed and Geometrically Distorted Images using Convolutional Neural Network', Journal of Computer Science, vol. 18, no. 8, 2022, doi: 10.3844/jcssp.2022.757.769.

[15] S. AbuRass, A. Huneiti, and M. B. Al-Zoubi, 'Enhancing Convolutional Neural Network using Hu's Moments', International Journal of Advanced Computer Science and Applications, vol. 11, no. 12, 2020, doi: 10.14569/IJACSA.2020.0111216.

[16] M. Alshawabkeh, M. H. Ryalat, O. M. Dorgham, K. Alkharabsheh, M. H. Btoush, and M. Alazab, 'A hybrid convolutional neural network model for detection of diabetic retinopathy', International Journal of Computer Applications in Technology, vol. 70, no. 3/4, p. 179, 2022, doi: 10.1504/IJCAT.2022.130886.

[17] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, 'ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators', Mar. 2020, [Online]. Available: http://arxiv.org/abs/2003.10555.

[18] D. Zhang and W. S. Lee, 'Question classification using support vector machines', in Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval, New York, NY, USA: ACM, Jul. 2003, pp. 26–32. doi: 10.1145/860435.860443.

[19] P. Zhou et al., 'Attention-Based Bidirectional Long Short-Term Memory Networks for Relation Classification', in Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), Stroudsburg, PA, USA: Association for Computational Linguistics, 2016, pp. 207–212. doi: 10.18653/v1/P16-2034.

[20] N. Reimers and I. Gurevych, 'Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks', in Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), Stroudsburg, PA, USA: Association for Computational Linguistics, 2019, pp. 3980–3990. doi: 10.18653/v1/D19-1410.

[21] W. Liu, P. Zhou, Z. Wang, Z. Zhao, H. Deng, and Q. JU, 'FastBERT: a Self-distilling BERT with Adaptive Inference Time', in Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, Stroudsburg, PA, USA: Association for Computational Linguistics, 2020, pp. 6035–6044. doi: 10.18653/v1/2020.acl-main.537.

[22] D. Britz, A. Goldie, M.-T. Luong, and Q. Le, 'Massive Exploration of Neural Machine Translation Architectures', in Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, Stroudsburg, PA, USA: Association for Computational Linguistics, 2017, pp. 1442–1451. doi: 10.18653/v1/D17-1151.

[23] J. Pennington, R. Socher, and C. D. Manning, 'GloVe: Global Vectors for Word Representation'. [Online]. Available: http://nlp.

[24] N. Van-Tu and L. Anh-Cuong, 'Improving Question Classification by Feature Extraction and Selection', Indian J Sci Technol, vol. 9, no. 17, May 2016, doi: 10.17485/ijst/2016/v9i17/93160.

[25] S. Chotirat and P. Meesad, 'Part-of-Speech tagging enhancement to natural language processing for Thai wh-question classification with deep learning', Heliyon, vol. 7, no. 10, p. e08216, Oct. 2021, doi: 10.1016/j.heliyon.2021.e08216.

[26] H. Tayyar Madabushi, M. Lee, and J. Barnden, 'Integrating Question Classification and Deep Learning for improved Answer Selection', in Proceedings of the 27th International Conference on Computational Linguistics, E. M. Bender, L. Derczynski, and P. Isabelle, Eds., Santa Fe, New Mexico, USA: Association for Computational Linguistics, Aug. 2018, pp. 3283–3294. [Online]. Available: https://aclanthology.org/C18-1278.

[27] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, 'Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling', Dec. 2014, [Online]. Available: http://arxiv.org/abs/1412.3555.

[28] A. Graves, 'Long Short-Term Memory', 2012, pp. 37–45. doi: 10.1007/978-3-642-24797-2_4.

[29] O. Sagi and L. Rokach, 'Ensemble learning: A survey', WIREs Data Mining and Knowledge Discovery, vol. 8, no. 4, Jul. 2018, doi: 10.1002/widm.1249.

[30] Z.-H. Zhou, Ensemble Methods Foundations and Algorithms. 2012.

[31] X. Dong, Z. Yu, W. Cao, Y. Shi, and Q. Ma, 'A survey on ensemble learning', Front Comput Sci, vol. 14, no. 2, pp. 241–258, Apr. 2020, doi: 10.1007/s11704-019-8208-z.

[32] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, 'BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding', Oct. 2018, [Online]. Available: http://arxiv.org/abs/1810.04805.

[33] Y. Liu et al., 'RoBERTa: A Robustly Optimized BERT Pretraining Approach', Jul. 2019, [Online]. Available: http://arxiv.org/abs/1907.11692.

[34] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, 'DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter', Oct. 2019, [Online]. Available: http://arxiv.org/abs/1910.01108.

[35] S. Aburass, 'Quantifying Overfitting: Introducing the Overfitting Index', 2023. Accessed: Nov. 10, 2023. [Online]. Available: https://arxiv.org/abs/2308.08682.

# SpanBERT-based Multilayer Fusion Model for Extractive Reading Comprehension

Pu Zhang, Lei He, Deng Xi

School of Computer Science and Technology
Chongqing University of Posts and Telecommunications, Chongqing, P. R. China

*Abstract*—**Extractive reading comprehension is a prominent research topic in machine reading comprehension, which aims to predict the correct answer from the given context. Pre-trained models have recently shown considerable effectiveness in this area. However, during the training process, most existing models face the problem of semantic information loss. To address this problem, this paper proposes a model based on the SpanBERT pre-trained model to predict answers using a multi-layer fusion method. Both the outputs of the intermediate layer and the prediction layer of the transformer are fused to perform answer prediction, thereby improving the model's performance. The proposed model achieves F1 scores of 92.54%, 84.02%, 80.86%, 71.32%, and EM scores of 86.27%, 81.25%, 69.10%, 56.42% on the SQuAD1.1, SQuAD2.0, Natural Questions and NewsQA datasets, respectively. Experimental results show that our model outperforms a number of existing models and has excellent performance.**

*Keywords—Machine reading comprehension; pre-trained model; transformer*

## I. INTRODUCTION

With the advent of the big data era, getting the right answers from massive amounts of data in a timely and accurate manner has become an urgent task. As one of the most popular natural language processing (NLP) tasks in recent years, machine reading comprehension (MRC) aims to enable machines to learn how to read and understand texts, that is, to find answers to given questions from relevant articles. Extractive reading comprehension, a subtask of MRC, has also made significant progress in recent years, requiring models to extract a continuous passage of text from the given input text as the final answer.

Extractive reading comprehension can overcome the limitations of relying solely on individual words or entities to answer questions. Its task is to use a model to extract an answer from a given passage or paragraph based on a given question. As shown in Fig. 1, given the question "when does season 2 of lethal weapon come out" and the passage "Lethal Weapon is an American buddy cop action comedy ...", the model reads and understands the passage and question, and then extracts a continuous segment "September 26, 2017" from the passage as the answer. The commonly used datasets for extractive reading comprehension include the SQuAD dataset [1], the NewsQA dataset [2], the TriviaQA dataset [3], and so on.

In deep learning-based reading comprehension models, the prediction of answer boundaries relies heavily on information interactions, and scholars have proposed one-way attention

models to employ attention mechanisms to enhance the interaction of information between passage and question. For example, Hermann et al. [4] used the one-way attention model as the basis for contextualizing questions, calculating the weight of each word in a passage, and generating a final representation of the passage. However, because one-way attention mechanisms only account for unidirectional attention, resulting in limited interaction between passage and question, researchers later proposed bidirectional attention models. For example, Seo et al. [5] proposed the bidirectional attention model BiDAF, which computes attention between question and passage separately in both directions to enhance information interaction and achieve good results.

In 2018, Google introduced the bidirectional pre-training language model BERT [6]. Since its release, BERT has achieved outstanding results in the field of NLP. Its remarkable performance is attributed to its internal multi-layer transformer structure, and directly using BERT with a simple answer predictor can achieve good performance on the SQuAD dataset [6].

The study by Ganesh et al. [7] showed that different layers of the transformer encoder focus on different semantic information. Ramnath S et al. [8] experimentally demonstrated that the prediction layer of BERT focuses more on contextual understanding and answer prediction, but ignores the interaction between context and question, while the earlier layers of the transformer focus more on the latter. Thus, some semantic information is lost during the learning process.

Since BERT was proposed, a number of pre-training models have been successively proposed. Among them, SpanBERT [9] is a representative model. Compared to BERT, SpanBERT [9] is more suitable for extractive tasks. However, it does not take into account the representational information emphasized by intermediate transformer layers, which may affect the final answer extraction in subsequent iterations.

This paper aims to address the potential problem of semantic information loss during learning in SpanBERT by investigating the fusion of semantic information from the intermediate and prediction layers of the transformer. A method is proposed to predict answers using a multi-layer transformer based on the SpanBERT model. The lower layer and the prediction layer of the transformer work together to predict answers. The predicted answer span information from both layers is combined to improve the accuracy of the predicted answers. The contributions of this work are outlined as follows:

*1)* We propose a model that can address the problem of semantic information loss during the learning process of the pre-trained model SpanBERT. Our model enhances the interaction between the paragraph text and the question by utilizing the SpanBERT model, the representation information emphasized by the intermediate layer and the final prediction layer of the transformer can be fused to improve the performance of the model's answer extraction.

*2)* A new approach to vector fusion is proposed in this study, which can effectively combine semantic information. An attention mechanism is utilized to fuse the outputs of the intermediate and prediction layers of the transformer, resulting in a new fused vector. This vector can be used to generate a probability distribution vector of the answer span, which is then multiplied by the answer prediction vector to obtain the predicted answer span. Combining the representational information that the final prediction layer and the middle layer focused on will improve the accuracy of the model's answer extraction.

*3)* We conduct comparison experiments on four datasets, including SQuAD1.1, SQuAD2.0, NaturalQA, and NewsQA, with two evaluation metrics, including F1 score and EM score, experiment results show that the proposed model has excellent performance.

| Passage | Lethal Weapon is an American buddy cop action comedy - drama television series that is based on the film series of the same name created by Shane Black . The series was ordered on May 10 , 2016 and premiered on Fox on September 21 , 2016 . On October 12 , 2016 , Fox picked up the series for a full season of 18 episodes . On February 22 , 2017 , Fox renewed the series for a 22 - episode second season , which premiered on September 26 , 2017. |
|---|---|
| Quesiton | when does season 2 of lethal weapon come out? |
| Answer | September 26 , 2017 |

Fig. 1. Example of extractive reading comprehension.

## II. RELATED WORK

MRC technology was first developed in the 1970s. In 1977, W.G. Lehnert et al. [10] designed the QUALM system, which used question-answering rules. In the 21st century, researchers integrated machine learning methods into MRC research. However, there were still drawbacks such as weak model generalization and insufficient feature extraction. The development of neural networks provided an opportunity for the advancement of MRC technology. From unidirectional attention mechanisms to bidirectional attention mechanisms, MRC technology has made significant strides and remarkable advancements. In recent years, the bidirectional pre-training language model BERT has achieved superior results in multiple task domains.

Reading comprehension tasks can be categorized as cloze tests, multiple-choice, span extraction, and generative reading comprehension.

Cloze-style tests prompt the machine to select the correct answer from a finite number of alternatives by removing words from the sentence. Representative datasets include CBT [11] and CNN/Daily Mail [4]. Representative models include the Gated Attention Reader [12] and others.

Multiple-choice reading comprehension has a more flexible answer format than cloze tests, as it is not limited to words or entities in context. However, the answers to the questions must still be provided in advance. Representative datasets for this type of task include MCTest [13] and RACE [14], while representative models include DCMN+ [15].

Currently, span extractive reading comprehension is the most popular task in this field, which is more challenging than traditional machine reading comprehension. The goal is to extract a contiguous span from a given text paragraph, which is not selected from a list of options.

Extractive reading comprehension models typically consist of four network architecture components: an embedding module, a feature extraction module, an information interaction module, and an answer prediction module. The embedding module converts each word in the passage and question into a fixed-length vector representation. To achieve this, a classical word vector encoding method such as Word2vec [16] can be used. The feature extraction module is often positioned after the embedding layer to extract context and question features separately. This module typically uses classical deep neural networks, such as recurrent neural networks and convolutional neural networks, to extract contextual information. The information interaction module is responsible for combining the encoded information of the paragraph and the question. It also captures the relationships between the words in the paragraph and the question to obtain their representations. The answer prediction module is located at the end of MRC systems and provides answers to questions based on the primary context.

In the extractive reading comprehension task, two classifiers are typically trained to predict the starting and ending indices of the answer. Common datasets for extractive reading comprehension include SQuAD, NewsQA, TriviaQA, SearchQA [17], and so on. Representative models include SpanBERT, BLANC [18], etc.

Although extractive reading comprehension has made significant advancements, its capabilities remain insufficient. Specifically, confining answers to a specific span within the context is still unrealistic. Generative reading comprehension requires machines to infer, summarize and provide open-ended answers from multiple passages of text. Of the four different types of tasks, generative reading comprehension is the most difficult. NarrativeQA [19] is a dataset that represents generative reading comprehension, and UniLMv2 [20] is a model that represents this type of task.

In addition to the task form, reading comprehension models can be structurally divided into a reading comprehension module and an answer prediction module. The reading comprehension module aims to answer the given questions based on the given passages. It is considered the core part of the model, where the model learns information from the input

text passage and question and generates the input text representation. For instance, Seo et al. [5] proposed BiDAF, which uses a bi-directional attention mechanism to improve the interaction between the question and the text passage, resulting in a more effective representation of the input text. BERT, on the other hand, enhances word embedding through multiple layers of transformer. SpanBERT, which is built on top of the BERT architecture, further improves text comprehension in the continuous span extraction task by training with span mask and span boundary objective. In this paper, we use SpanBERT as the foundational architecture of our model.

The answer prediction module is divided into different types of tasks. For the cloze reading comprehension task, the module predicts the probability values of multiple candidate answers based on the text vector information and selects the option with the highest probability as the predicted answer. For the extractive reading comprehension task, the answer is a continuous segment of the given text. The answer prediction module should generate two probability distributions based on the text representation: one for the starting position of the answer and the other for the ending position.

### III. MODEL

Our model utilizes SpanBERT, a pre-trained model with 12 layers of transformer encoder, similar to BERT. The 12th layer is typically used for final answer prediction. However, language is complex and contains not only grammar and semantic information, but also hidden information such as emotion and deduction. Therefore, each layer of the transformer encoder learns different information during the training process. This model addresses the limitation of existing models that ignore other potentially helpful layers for answer prediction, by incorporating them into the prediction process. To improve answer prediction performance, this paper proposes a model that utilizes an information fusion approach. The basic architecture of our model is shown in Fig. 2.

First, both the passage and the question are input into the embedding layer for learning, resulting in the output of each layer of the transformer. Subsequently, the output of the middle layer and the output of the prediction layer are combined to obtain a fusion vector that contains the semantic information from both the middle layer and the prediction layer. Finally, the answer interval information predicted from the fusion vector is further fused with the answer information extracted from the prediction layer to obtain the final answer.

A thorough explanation of our model, including the encoding layer, encoder attention block, answer extraction, loss function, and other components, is given in this section.

### A. Embedding

Suppose the question sequence is $Q=[q_1, q_2, q_3, \ldots, q_m]$ and the passage sequence is $P=[p_1, p_2, p_3, \ldots, p_n]$. They are separated by a separator when inputting to the model, as shown in the following formula:

$$[\text{CLS}], q_1, q_2, q_3, \ldots, q_m, [\text{SEP}], p_1, p_2, p_3, \ldots, p_n, [\text{SEP}] \# \quad (1)$$

After inputting the text into SpanBERT, the encoding sequences $T_n$ ($n = 1\sim12$) of each layer of the transformer encoder can be obtained through the encoder modules. When $n$

$= 12$, $T_{12}$ represents the encoding sequence of the SpanBERT prediction layer, and the example of $T_n$ is as follows:

$$T_n = T_{[\text{CLS}]}, T_{q_1}, T_{q_2}, \ldots, T_{q_m}, T_{[\text{SEP}]}, T_{p_1}, T_{p_2}, \ldots, T_{p_n}, T_{[\text{SEP}]} \quad (2)$$



Fig. 2. Model architecture.

### B. Encoder Attention

The specific steps of information fusion are as follows:

*1)* Take out the results of the n-th layer and the last layer of the transformer encoder, and the encoding of the question is separately extracted to obtain the encoding information of the question for the two layers, denoted as $Q_n$ and $Q_{Last}$, respectively.

*2)* The semantic information is combined between the two layers by means of dot product, and then calculate the weight of each of the two layers through the fully connected layer. The formula is as follows:

$$Q'_n = SUM(Q_n) \quad (3)$$

$$Q'_{Last} = SUM(Q_{Last}) \quad (4)$$

$$E_{Attention} = Q'_n * Q'_{Last} \quad (5)$$

$$w_n, w_{Last} = Linear(E_{Attention}) \quad (6)$$

$Q_n$ and $Q_{Last}$ represent the encoded vectors of the question part in the *n*-th layer and the prediction layer of the transformer encoder. $SUM(*)$ is used to avoid the problem of inconsistent lengths of $Q$ in the input text by stacking the encoding information of $Q_n$ and $Q_{Last}$ according to the word encoding dimension. $E_{Attention}$ represents the fused semantic information. $Linear(*)$ is a linear function. $w_n$ and $w_{Last}$ represent the weights calculated for the *n*-th layer and prediction layer when predicting the answer, respectively.

*3)* By calculating the weights, we can obtain the fusion vector encoding as follows.

$$T_{merge} = w_n * T_n + w_{Last} * T_{Last} \qquad (7)$$

$T_n$ and $T_{Last}$ represent the encoding sequences obtained from the *n*-th layer and the *12*-th layer of the transformer encoder, respectively. $T_{merge}$ represents the fused vector of the word encoding information.

*C. Answer-span Prediction*

In the previous section, we fused the encoding of the *n*-th layer and the prediction layer to obtain a new fused vector. In this section, we propose a new approach for span prediction to improve the performance of the model. The formula is as follows:

$$pred_s = \frac{\exp(W_a T_{merge} + b_s^a)}{\sum_j \exp(W_a T_{merge_j} + b_s^a)} \qquad (8)$$

$$pred_e = \frac{\exp(V_a T_{merge} + b_e^a)}{\sum_j \exp(V_a T_{merge_j} + b_e^a)} \qquad (9)$$

$T_{merge}$ represents the fusion vector that fuses the word encoding information from the *n*-th and prediction layers. $W_a$、$V_a$、$b_s^a$ and $b_e^a$ represent the trainable weights and bias parameters. $pred_s$ and $pred_e$ represent the relative probability of each word as the beginning and ending position of the answer.

After obtaining the probability distributions of the answer span, which are represented by the fused vector of $pred_s$ and $pred_e$, the calculation of the guided-layer attention vector can be initiated. The formula is as follows:

$$Attention = pred_s * pred_e \qquad (10)$$



Fig. 3. Distribution of guided-layer attention vector.

The guided-layer vector represents the probability distribution of the predicted answer region. By multiplying $pred_s$ and $pred_e$, as shown in Fig. 3, the predicted probability value of words belonging to the answer span is higher than both ends, and the farther the distance is, the lower the probability value is. That is, the overall distribution of *Attention* is normal. The overall distribution exhibits a normal distribution, which could facilitate the ability to predict answer span for the guided-layer attention vector.

Finally, the answer prediction layer $T_{Last}$ and the probability distribution vector of the answer region are dot-multiplied, and then the final answer prediction is calculated by $softmax(*)$. The calculation formulas are as follows.

$$logits_s, logits_e = Split(T_{Last}) \qquad (11)$$

$$Ans_s = logits_s * Attention \qquad (12)$$

$$Ans_e = logits_e * Attention \qquad (13)$$

$$p_{i=s_a} = \frac{\exp(W_a Ans_{s_i} + b_s^a)}{\sum_j \exp(W_a Ans_{s_j} + b_s^a)} \qquad (14)$$

$$p_{i=e_a} = \frac{\exp(V_a Ans_{e_i} + b_e^a)}{\sum_j \exp(V_a Ans_{e_j} + b_e^a)} \qquad (15)$$

$Split(*)$ represents the vector splitting operation. $logits_s$ and $logits_e$ are used to predict the start and end indices of context. $Ans_s$ and $Ans_e$ denote answer prediction vectors that have merged the information of the probability distribution vector of answer region. $W_a$、$V_a$、$b_s^a$ and $b_e^a$ represent the trainable weights and bias parameters.

*D. Loss Function*

The loss function used in our model is the joint cross-entropy loss function, which is based on the negative log probabilities of the true answer's start and end positions in the predicted distribution. The formula is as follows:

$$L_{12} = -\frac{1}{N} \sum_{i=1}^{N} \left[ \log\left(P_{y_i^s Last}^{start}\right) + \log\left(P_{y_i^e Last}^{end}\right) \right] \qquad (16)$$

$P^{start}$ and $P^{end}$ are the probability distributions of the start and end positions of the answers predicted by the model. *Last* represents the prediction layer. $y_i^s$ and $y_i^e$ are the start and end positions of the real answer in the *i*-th training sample.

The same answer prediction is done for the fusion vector to obtain the loss function $L_n$.

$$L_n = -\frac{1}{N} \sum_{i=1}^{N} \left[ \log\left(P_{y_{i_n}^s}^{start}\right) + \log\left(P_{y_{i_n}^e}^{end}\right) \right] \qquad (17)$$

We define our final loss function as the weighted sum of the two loss functions:

$$L_{total} = (1 - \lambda)L_{12} + \lambda L_n \qquad (18)$$

$\lambda$ is a hyper-parameter moderating the ratio of two loss functions.

## IV. EXPERIMENTAL SETUP

*A. Datasets*

To validate the effectiveness of the proposed model, experiments were conducted and analyzed on four datasets: SQuAD1.1, SQuAD2.0 [21], Natural Questions [22] and NewsQA.

The Stanford Question Answering Dataset (SQuAD) is a large-scale English reading comprehension dataset constructed by Stanford University, which has been an indispensable dataset for MRC tasks since its release and has a milestone significance for the development of MRC technology. SQuAD1.1 contains 536 high-quality articles from English Wikipedia, which are divided into natural paragraphs. In

addition, there are 107,785 questions and corresponding answers, all of which are manually annotated.

SQuAD 2.0 builds on SQuAD 1.1 by adding unanswerable questions. Dataset creators provide an unanswerable question for each paragraph to interfere with the model's prediction. The training dataset contains 87k answerable and 43k unanswerable questions.

Natural Questions is a dataset of natural language queries. Each example consists of a Google query and a Wikipedia passage, where the answer is a span of the Wikipedia passage. The Natural Questions dataset contains 300,000 natural questions with human annotated answers.

The NewsQA dataset contains examples selected from over 10,000 news articles from CNN, along with 119,633 manually generated question-answer pairs. The answers are snippets of any length from the news article, and the dataset also includes partially unanswerable questions.

### B. Evaluation Metrics

F1 is the most widely used evaluation metric in existing extractive reading comprehension models. The Precision value is computed as follows:

$$Precision = \frac{TP}{TP + FP} \tag{19}$$

*TP* represents the number of true positive samples and *FP* represents the number of false positive samples. Then the Recall value is then calculated as follows:

$$Recall = \frac{TP}{TP + FN} \tag{20}$$

*FN* represents the number of false negative samples. F1 value is calculated from Precision and Recall with the following formula.

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \tag{21}$$

EM (Exact Match) is a common evaluation metric for question answering systems, and it is also one of the main metrics for SQuAD. It measures the percentage of all predictions that exactly match the ground-truth answer. The calculation formula is as follows:

$$EM = \frac{N_{right}}{N_{all}} \tag{22}$$

$N_{right}$ indicates the number of correct predictions and $N_{all}$ represents the number of all predictions.

### C. Experimental Setup

The implementation of this model is based on Python and its third-party libraries, with PyTorch serving as the deep learning framework. The hyperparameters λ are set to 0.8 in joint loss functions. The experimental datasets used are SQUAD1.1, SQUAD2.0, NaturalQA, and NewsQA. To facilitate model performance testing and due to limited computing resources, the training batch size is set to 16 for the SQUAD2.0 dataset and 8 for the other three datasets. The learning rate is set to $2*e^{-5}$, and the maximum length of input text is set to 384. The word vector dimension is set to 768, and the number of training epochs is set to 4 on the SQuAD2.0 dataset and 3 on the other three datasets.

### D. Baselines

BLANC [18]: To improve the accuracy of the final answer prediction, the model primarily employs a context prediction method. The model first predicts a soft label, then uses this soft label to calculate the context boundary probability, and finally uses the context boundary probability to optimize the final answer boundary prediction.

BERT-base [6]: BERT-base is a pre-trained language representation model that generates a deep bidirectional language representation using a masked language model (MLM). It is regarded as a landmark model in MRC, significantly advancing the field's development.

SpanBERT [9]: This BERT variation is tuned for fragment extraction tasks, resulting in more accurate representations. Two aspects contribute to the optimization. Firstly, it recommends adopting span masking rather than single-word masking for learning at fragments. Second, it trains the masked boundary words representation to anticipate masked fragment information.

ALBERT [23]: Compared to BERT, ALBERT overcomes the difficulties of extensive model parameterization and growing training time. It incorporates three major innovations: factorized embedding parameters, shared parameters across layers, and Sentence Order Prediction (SOP). The SOP creates not only positive examples by establishing the correct order of two consecutive sentences, but also negative examples by reversing their order.

LinkBERT [24]: LinkBERT is a cross-document language modeling training method that takes advantage of document links. In contrast with BERT, this approach has a distinct benefit in that it uses the links between documents to improve language modeling. LinkBERT treats the corpus as a document graph and employs linked documents as supplementary input to the model rather than modeling a single document.

DeBERT-base [25]: The proposed model aims to increase the robustness and effectiveness of the system when dealing with incomplete data. This is achieved by reconstructing hidden embeddings for sentences containing missing words.

OneS [26]: This model is based on the human learning model and introduces a new task of extracting essential knowledge from different knowledge sets for model pre-training.

KALA [27]: To solve the problem of catastrophic forgetting that occurs during adaptive pre-training, this model adjusts the intermediate hidden layer representation of a pre-trained model by incorporating knowledge from multiple domains.

ALBERTbase+V4ES [28]: This model proposes a verification mechanism that divides the machine learning process into two modules: general reading and fine-grained reading. The general reading module involves reading the text and question to obtain a preliminary answer. The fine-grained reading module reads again and generates a final answer.

RoBERTa-base [29]: This model enhances its performance by increasing the number of parameters and training data.

## V. RESULTS AND DISCUSSION

### A. Results

In this section, two evaluation metrics (F1 and EM) are used on four datasets to verify the effectiveness of the proposed model in the paper. The experimental results are shown in Table I, Table II, Table III, and Table IV.

From the experimental results on the SQuAD1.1 dataset, the results in Table I show that the proposed model achieves good performance. The F1 score of our model improved by 0.6% and the EM score improved by 0.84% compared to SpanBERT. Compared to other models such as BERT-base, ALBERT-large, DeBERT-base, OneS and BLANC, the F1 scores are improved by 4.04%, 1.94%, 0.44%, 2.84% and 0.67%, and the EM scores are improved by 5.47%, 2.37%, 0.17%, 3.27% and 0.97%, respectively. In the experimental results of the SQuAD2.0 dataset, as shown in Table II, the F1 score of our model is improved by 0.59% and the EM score is improved by 0.76% compared to SpanBERT. In addition, compared to models such as BERT-base, OneS and ALBERTbase+V4ES, the F1 scores of our model are improved by 3.62%, 3.82% and 0.62%, and the EM scores are improved by 3.65%, 4.15% and 0.95%, respectively.

From the results of the NaturalQA dataset, as shown in Table III, the F1 score of our model is improved by 2.55% and the EM score is improved by 2.50% compared to SpanBERT. Compared to other models such as BERT-base, ALBERT and BLANC, the F1 scores of the model can be improved by 4.47%, 4.97% and 0.82% respectively, and the EM scores can be improved by 4.62%, 5.29% and 0.77%, respectively. From the experimental results of the NewsQA dataset, as shown in Table IV, the F1 score of our model improved from 67.93% of SpanBERT to 71.32%, with an improvement of 3.39%, and the EM score of the model is improved by 3.57%.

All of the above results show that our model has excellent performance.

Overall, our model performs better on the NaturalQA and NewsQA datasets when compared to the SQuAD1.1 and SQuAD2.0 datasets. This is because the SQuAD datasets have a flaw where the questions and answers are very similar, resulting in the front layer of the transformer losing less semantic information during the learning process. On the other hand, the NaturalQA and NewsQA datasets are generated based on real questions and answers, so more semantic information is lost in the front layer transformer during model learning.

### B. Effect of Different Layers on Results

There are 12 layers of transformer encoders in SpanBERT, but each layer focuses on different information. During its iterative learning process, some information about the answer in the earlier transformer encoder layers may be forgotten. In this section, F1 and EM are used as evaluation metrics to investigate the effectiveness of each layer in guiding the prediction layer to predict the answer, using SQuAD1.1,

SQuAD 2.0, NaturalQA and NewsQA datasets as experimental datasets. The experimental results are shown in Fig. 4.

TABLE I.     RESULTS (%) OF EXPERIMENTS ON THE SQUAD1.1

| Models | F1 | EM |
|---|---|---|
| BERT-base | 88.50 | 80.80 |
| SpanBERT | 91.94 | 85.43 |
| ALBERT-large | 90.60 | 83.90 |
| LinkBERT | 90.10 | - |
| DeBERT-base | 92.10 | 86.10 |
| OneS | 89.70 | 83.00 |
| BLANC | 91.87 | 85.30 |
| ALBERTbase+V4ES | 91.10 | 83.40 |
| Our model | **92.54** | **86.27** |

TABLE II.     RESULTS (%) OF EXPERIMENTS ON THE SQUAD2.0

| Models | F1 | EM |
|---|---|---|
| BERT-base | 80.40 | 77.60 |
| SpanBERT | 83.43 | 80.49 |
| ALBERT-large | 82.30 | 79.40 |
| DeBERT-base | 82.50 | 79.30 |
| OneS | 80.20 | 77.10 |
| ALBERTbase+V4ES | 83.40 | 80.30 |
| RoBERTa-base | 83.70 | 80.50 |
| Our model | **84.02** | **81.25** |

TABLE III.     RESULTS (%) OF EXPERIMENTS ON THE NATURALQA

| Models | F1 | EM |
|---|---|---|
| BERT-base | 76.39 | 64.48 |
| SpanBERT | 78.31 | 66.60 |
| ALBERT-large | 75.89 | 63.81 |
| LinkBERT | 78.30 | - |
| BLANC | 80.04 | 68.33 |
| Our model | **80.86** | **69.10** |

TABLE IV.     RESULTS (%) OF EXPERIMENTS ON THE NEWSQA

| Models | F1 | EM |
|---|---|---|
| BERT-base | 65.07 | 50.11 |
| SpanBERT | 67.93 | 52.85 |
| ALBERT-large | 66.02 | 51.18 |
| LinkBERT | 69.30 | - |
| KALA | 68.27 | 54.25 |
| BLANC | 70.31 | 55.52 |
| Our model | **71.32** | **56.42** |

Fig. 4.   (a), (b), (c), (d) represent the experimental results (F1/EM score) of different layers of SpanBERT as a guided-layer on SQuAD1.1, SQuAD2.0, NaturalQA, and NewsQA datasets (%), respectively.

From Fig. 4(a), we can see that on the SQuAD1.1 dataset, our model achieves optimal performance when using the output of the 11-th layer of the transformer encoder for fusion vector calculation. Specifically, compared to the SpanBERT model, the F1 score increased by 0.6% to 92.54% and the EM score increased by 0.84% to 86.27%. From Fig. 4(b), it can be observed that when using the prediction layer of the SpanBERT to directly predict the answer, the model achieves 78.31% (F1) and 66.6%(EM) on the NaturalQA dataset. However, when using different intermediate layers of the SpanBERT with the prediction layer to generate a new fused vector through encoder attention, the model achieves an improvement of 2.04% to 2.41% in F1 score and 1.93% to 2.5% in EM score. Among these experiments, the best performance was achieved when using the 5-th layer. This suggests that there is some semantic information loss in the early layers of the SpanBERT encoder during the iterative process. It also demonstrates the effectiveness of the proposed model.

As shown in Fig. 4(c) and Fig. 4(d), experiments on the SQuAD2.0 and NewsQA datasets, it was found that the best performance could be achieved when using the outputs of the 6-th and 2-nd layers of transformer encoder, respectively. Choosing different layers of transformer encoder to conduct experiments on different datasets always resulted in the best performance, indicating that our model is effective, and the 12 layers of transformer encoder in the SpanBERT model have differences in processing semantic information, and each layer focuses on different semantic information.

### C. Effects of Hyperparameter $\lambda$

In the experiment, a weighted sum of the joint loss function $L_{merge}$ and the loss function $L_{12}$ of the answer prediction layer

was used as the total loss function of the model, and a hyperparameter $\lambda$ was introduced to balance the weight of the two loss values. In this section, we use SQuAD1.1 and NewsQA as experimental datasets, with F1 score and EM score as evaluation metrics, to verify the optimal value of the hyperparameter $\lambda$ in the joint loss function. The experimental results are shown in Fig. 5.

We set $\lambda$ to [0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0], and conduct experiments by incorporating $\lambda$ into the calculation of the joint loss function $L_{total}$. As $\lambda$ increases, the accuracy of the model in predicting answers increases. In the experiments on the SQuAD1.1 dataset, the model achieved the best performance when $\lambda$ was set to 0.8, with an F1 score of 92.54% and an EM score of 86.27%. In the experiments on the NewsQA dataset, the model achieved the best performance when $\lambda$ was set to 0.7, with an F1 score of 71.35% and an EM score of 56.13%. As $\lambda$ increases further, the model's performance decreases.



Fig. 5.   (a) Results (%) for different values of $\lambda$ on SQuAD1.1 dataset. (b) Results (%) for different values of $\lambda$ on NewsQA dataset.

## D. Different Pre-trained Models

In the above experiments, we used SpanBERT as a pre-training model to verify the effectiveness of our model. In order to verify the applicability of this method to other pre-trained models, we also conducted comparative experiments by using two pre-trained models including BERT and SpanBERT. The experimental results are shown in Table V and Table VI.

TABLE V.     RESULTS WITH DIFFERENT PRE-TRAINED MODELS (SQUAD1.1)

| | | SQuAD1.1 | |
|---|---|---|---|
| | | F1 | EM |
| Bert-BASE | BASELINE | 88.10 | 80.49 |
| | +our method | **88.76** | **81.34** |
| SpanBERT | BASELINE | 91.58 | 84.97 |
| | +our method | **92.54** | **86.27** |

TABLE VI.     RESULTS WITH DIFFERENT PRE-TRAINED MODELS (NEWSQA)

| | | NewsQA | |
|---|---|---|---|
| | | F1 | EM |
| Bert-BASE | BASELINE | 65.07 | 50.11 |
| | +our method | **66.74** | **51.69** |
| SpanBERT | BASELINE | 67.93 | 52.85 |
| | +our method | **71.32** | **56.42** |

According to Table V, we can see that when using BERT as the pre-training model, our model can improve the model's F1 score by 0.66% to reach 88.76%. The model's EM score can be improved by 0.85% to reach 81.34%. When using SpanBERT as the pre-training model, the model's F1 score can be improved by 0.96% to reach 92.54%, and the model's EM score can be improved by 1.30% to reach 86.27%. From Table VI, on the NewsQA dataset, when using BERT as the pre-training model, our model can improve the model's F1 score and EM score from 65.07% and 50.11% to 66.74% and 51.69%, respectively. Using SpanBERT as the pre-training model, the model's F1 score can be improved by 3.39% to reach 71.32%, and the model's EM score can be improved by 3.57% to reach 56.42%. These results suggest that our method is applicable to other pre-training models.

## E. Effects of $T_{merge}$

In this section, we conduct comparative experiments using just the fusion vector $T_{merge}$ as the prediction layer and evaluate the performance. The experimental results are shown in Table VII and Table VIII.

TABLE VII.     RESULTS OF $T_{merge}$ (NATURALQA)

| | NaturalQA | |
|---|---|---|
| | F1 | EM |
| $T_{merge}$ | 78.96 | 67.21 |
| SpanBERT | 78.31 | 66.60 |
| Our model | **80.86** | **69.10** |

TABLE VIII.     RESULTS OF $T_{merge}$ (NEWSQA)

| | NewsQA | |
|---|---|---|
| | F1 | EM |
| $T_{merge}$ | 68.54 | 53.40 |
| SpanBERT | 67.93 | 52.58 |
| Our model | **71.32** | **56.42** |

According to the experimental results shown in Table VII, on the NaturalQA dataset, when using the fusion vector $T_{merge}$ as the prediction layer, the model achieved an F1 score of 78.96% and an EM score of 67.21%, which is higher than the pre-trained model SpanBERT's 78.31% and 66.60%. However, there is still a performance gap compared to the experimental results of the proposed model in this paper. Similar experimental results were also verified on the NewsQA dataset. We can see that our model still has the highest performance, which indicates that our model is effective.

## VI.   CONCLUSION

In this paper, a SpanBERT-based multi-layer fusion extractive reading comprehension model is proposed. By fusing the representational information obtained from the intermediate transformer layer with the representational information obtained from the prediction layer, a new fusion vector is obtained through an encoder attention mechanism. Using the fusion vector, the distribution probability vector of the answer region is then computed and used together with the prediction layer to jointly predict the answer. Finally, answer extraction is performed. We have conducted extensive experiments to demonstrate the effectiveness of our model. Although the comparative experiments have shown the clear performance advantages of our model on the respective datasets, there are still certain problems and room for improvement. Specifically, even though the learning process of the pre-training model no longer suffers from the loss of semantic information, the learning process of our model still relies solely on the input data and does not use external knowledge, whereas people often use external knowledge to improve their understanding of textual data during the reading comprehension process. As a result, future studies can use the incorporation of external knowledge to enhance the semantic data, thereby improving the performance of the model. In the future work, we will also explore other pre-trained models for machine reading comprehension tasks.

## REFERENCES

[1]    Rajpurkar P, Zhang J, Lopyrev K, Liang P. (2016). SQuAD: 100,000+ Questions for Machine Comprehension of Text[C]. Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, 2383-2392.

[2]    Adam Trischler, Tong Wang, Xingdi Yuan, et al. (2017). NewsQA: A Machine Comprehension Dataset[C]. Proceedings of the 2nd Workshop on Representation Learning for NLP, 191-200.

[3]    Joshi M, Choi E, Weld DS, Zettlemoyer L. (2017). TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension[C]. Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, 1601-1611.

[4]    Hermann K M, Kocisky T, Grefenstette E, et al. (2015). Teaching machines to read and comprehend[C]. Proceedings of the 28th

International Conference on Neural Information Processing Systems, 1693-1701.

[5] Seo M., Kembhavi A., Farhadi A., et al. (2017). Bidirectional attention flow for machine comprehension[C]. Proceedings of the 5th International Conference on Learning Representations, 1437-1450.

[6] Devlin J, Chang M W, Lee K, Toutanova K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding[C]. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics, 4171- 4186.

[7] Ganesh J, Sagot B, Seddah D. (2017). What does BERT learn about the structure of language?[C]. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 3651–3657.

[8] Ramnath S, Nema P, Sahni D, et al. (2020).Towards interpreting BERT for reading comprehension based QA[C]. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 3236–3242.

[9] Joshi M., Chen D., Liu Y., et al. (2020). Spanbert: Improving pre-training by representing and predicting spans[J]. Transactions of the Association for Computational Linguistics, 64-77.

[10] Lehnert W. G. The process of question answering[M]. Yale University, 1977: 35-76.

[11] Hill F., Bordes A., Chopra S., et al. (2016). The Goldilocks Principle: Reading children's books with explicit memory representations[C]. Proceedings of the 4th International Conference on Learning Representations, 1124-1137.

[12] Dhingra B, Liu HX, Yang ZL, Cohen W, Salakhutdinov R. (2017). Gated-Attention Readers for Text Comprehension[C]. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, 1832–1846.

[13] Richardson M，Burges C J C，Renshaw E. (2013). MCTest：a challenge dataset for the open-domain machine comprehension of text[C]. roceedings of the 2013 Conference on Empirical Methods in Natural Language Processing，193-203.

[14] Lai GK, Xie QZ, Liu HX, Yang YM, and Hovy E. (2017). RACE: Large-scale ReAding Comprehension Dataset From Examinations [C]. Proceedings of the 2017 conference on empirical methods in natural language processing,785-794.

[15] Zhang S, Zhao H, Wu Y, et al. (2020). DCMN+: Dual co-matching network for multi-choice reading comprehension[C]. Proceedings of the AAAI Conference on Artificial Intelligence, 9563-9570.

[16] Mikolov T., Chen K., Corrado G., et al. (2013). Efficient estimation of word representations in vector space[C]. Proceedings of the 1st International Conference on Learning Representations, 976-988.

[17] Dunn M, Sagun L, Higgins M, et al. SearchQA: A new Q&A dataset augmented with context from a search engine[J]. arXiv preprint arXiv:1704.05179, 2017.

[18] Seonwoo Y., Kim J. H., Ha J. W.,Oh A. (2020). Context-Aware answer extraction in question answering[C]. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2418-2428.

[19] Kociský T, Schwarz J, Blunsom P, et al. (2018). The NarrativeQA Reading Comprehension Challenge [J]. Transactions of the Association for Computational Linguistics, 317-328.

[20] Bao H, Dong L, Wei F, et al. (2020). Unilmv2: Pseudo-masked language models for unified language model pre-training[C]. International Conference on Machine Learning, 642-652.

[21] Rajpurkar P., Jia R., Liang P. (2018). Know What You Don't Know: Unanswerable questions for SQuAD[C]. Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, 784-789.

[22] Kwiatkowski T, Palomaki J, Redfield O, et al. (2019). Natural Questions: a benchmark for question answering research[J]. Transactions of the Association for Computational Linguistics, 453-466.

[23] Lan Z., Chen M., Goodman S., et al. (2020). ALBERT: A Lite BERT for Self-supervised Learning of Language Representations[C]. Proceedings of the 8th International Conference on Learning Representations, 1362-1379.

[24] Yasunaga M, Leskovec J, Liang P. (2022). LinkBERT: Pretraining Language Models with Document Links[C]. Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics , 8003–8016.

[25] He P., Liu X., Gao J., et al. (2021). DeBERTa: Decoding-enhanced bert with disentangled attention[C]. Proceedings of the 9th International Conference on Learning Representations, 1278-1301.

[26] Xue FZ, He XX, Ren XZ, et al. One Student Knows All Experts Know: From sparse to dense[J]. arXiv preprint arXiv:2201.10890, 2022.

[27] Kang M, Baek J, Hwang S J. (2022). KALA: Knowledge-Augmented Language Model Adaptation[C]. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 5144–5167.

[28] Peng Y, Li XY, Song JK, et al. (2021). Verification mechanism to obtain an elaborate answer span in machine reading comprehension[J]. Neurocomputing, 80-91.

[29] Liu Y, Ott M, Goyal N, et al. Roberta: A robustly optimized bert pretraining approach[J]. arXiv preprint arXiv:1907.11692, 2019.

# Topology Approach for Crude Oil Price Forecasting of Particle Swarm Optimization and Long Short-Term Memory

Marina Yusoff[1], Darul Ehsan[2], Muhammad Yusof Sharif[3], Mohamad Taufik Mohd Sallehud-din[4]

Institute for Big Data Analytics and Artificial Intelligence (IBDAAI), Kompleks Al-Khawarizmi,
Universiti Teknologi MARA (UiTM), 40450 Shah Alam, Selangor[1]
College of Computing, Informatics and Mathematics, Kompleks Al-Khawarizmi,
Universiti Teknologi MARA (UiTM), 40450 Shah Alam, Selangor[2]
College of Computing, Informatics and Mathematics, Universiti Teknologi MARA Shah Alam, Malaysia[3]
PETRONAS Research Sdn Bhd, Jln Ayer Hitam, Kawasan Institusi Bangi, 43000 Bandar Baru Bangi, Selangor[4]

*Abstract*—Forecasting crude oil prices hold significant importance in finance, energy, and economics, given its extensive impact on worldwide markets and socio-economic equilibrium. Using Long Short-Term Memory (LSTM) neural networks has exhibited noteworthy achievements in time series forecasting, specifically in predicting crude oil prices. Nevertheless, LSTM models frequently depend on the manual adjustment of hyperparameters, a task that can be laborious and demanding. This study presents a novel methodology incorporating Particle Swarm Optimization (PSO) into LSTM networks to optimize the network architecture and minimize the error. This study employs historical data on crude oil prices to explore and identify optimal hyperparameters autonomously and embedded with the star and ring topology of PSO to address the local and global search capabilities. The findings demonstrate that LSTM+starPSO is superior to LSTM+ringPSO, previous hybrid LSTM-PSO, conventional LSTM networks, and statistical time series methods in its predictive accuracy. LSTM+starPSO model offers a better RMSE of about +0.16% and +22.82% for WTI and BRENT datasets, respectively. The results indicate that the LSTM model, when enhanced with PSO, demonstrates a better proficiency in capturing the patterns and inherent dynamics data changes of crude oil prices. The proposed model offers a dual benefit by alleviating the need for manual hyperparameter tuning and serving as a valuable resource for stakeholders in the energy and financial industries interested in obtaining dependable insights into fluctuations in crude oil prices.

*Keywords—Crude oil; deep learning; Particle Swarm Optimization; Long Term-Short Memory; forecasting*

## I. INTRODUCTION

As one of the most significant commodities in the world, crude oil is responsible for the energy consumption. It is the foundation for daily items, from plastics to transportation fuels. Considering that fluctuations in crude oil prices significantly influence economies worldwide, price forecasting can help reduce the risks of oil price volatility [1]. Predictive methods in oil and gas operations can boost efficiency, lower costs, and reduce environmental impact from a good forecasting model [2]. Machine learning researchers and developers face challenges when working with large datasets and diverse data types, primarily because of noisy and unclean data [3]. Several pre-processing methods have been

developed to address this issue, with specific methods yielding favourable outcomes. Hence, the choice of pre-processing techniques would depend upon the data's characteristics and quality. Typically, benchmark datasets such as ready data do not necessitate extensive pre-processing tasks [4]. However, the most significant challenge is effectively managing substantial quantities, especially the time series data, which requires the development of a more understandable model. The abovementioned difficulties are relevant to the oil and gas data, especially in real-time data monitoring. Further investigation is required to effectively tackle the obstacles and determine the practicability of these methodologies using benchmark and real-life data.

Recently, there has been an increasing preference for incorporating predictive analytics in the oil and gas sector. Machine learning methods and their diverse applications in the oil and gas sector encompass multiple areas, such as pipeline prediction [5], well-log formation [6], and crude oil price forecasting [7].

Due to its effectiveness in predictive analytics, the Long Short-Term Memory (LSTM) model is widely utilized in various oil and gas industry sub-fields and other engineering and finance-related disciplines. Prior studies have examined various methodologies on LSTM, including LSTM with optimization and CNN with LSTM. For instance, CNN and LSTM address two instances of degradation prediction in offshore operation platforms for natural gas treatment plants and seawater injection pumps for oil [8]. Compared to a single LSTM model, the performance of the CNN with the LSTM model is superior, exhibiting a notable enhancement of 15.5% in precision.

Furthermore, the performance of LSTM has also been documented in reputable studies on time series [9], [10]. Yang et al. [9] employed the LSTM model to forecast short-and long-term production events in shale gas wells. The method performed better than the ARIMA, Arps, and Duong methods. In another study, Song et al. [11] used LSTM, a feature extraction and optimization model that incorporates feature engineering and parameter optimization and exhibits the lowest mean absolute error (MAE) value compared to other

models such as BPNN, LSTM, and random forest as reported by Dyer et al. [12]

One of the recent challenges in oil and gas is predicting crude oil prices. The demand for crude oil price prediction has increased due to crude oil's complex and highly unpredictable characteristics of crude oil [13]. Several methods were proposed and evaluated using benchmark datasets. For instance, LSTM and Henry gas solubility optimization (CHGSO) technique to estimate crude oil prices using West Texas Intermediate (WTI) and Brent Crude Oil Time Series COTS datasets [14] and Hybrid Wavelet Transform (WT) Bidirectional Long Short-Term Memory Network (BiLSTM)-Attention-CNN. WT-BLAC performs well for WTI data with R2, RMSE, MAPE, and MAE of 0.97, 2.25, 1.18, and 2.63, respectively. Furthermore, it evaluated a similar dataset using ensemble and ANN but with a different range of time series data [15]. The models have acceptable and significant interpretability in time series prediction of crude oil futures prices.

Support vector regression model did another forecasting crude oil prices, which are infamous for being unpredictable and have been fine-tuned using a genetic algorithm [2]. They use a ten-year daily dataset from NASDAQ and key economic input features. A study by Shahbazbegian et al. [10] divided time series into sub-series by the proposed hybrid model, which employs a multifaceted approach to capture distinct characteristics, LSTM is combined with the Markov switching model to forecast volatile and fluctuating sub-series.

After integrating these predictions using a linear combination, a comprehensive estimation of the time series for WTI crude oil prices is generated. The proposed method's respective RMSE and MAPE values are 4.18 and 0.03. Furthermore, He et al. [16] found that a hybrid forecasting model based on multi-modal features for price trends and employing the variational mode decomposition algorithm, extraction of data features with multiple modes, and time series employment of analysis provides acceptable performance. More research on crude oil forecasting solutions is still in demand. LSTM and its variants have great potential to obtain better forecasting results by embedding an appropriate optimization method. This paper focuses on using LSTM embedded with one of the popular computational optimizations, PSO. PSO is chosen due to its ease of implementation, high precision, and fast convergence [17], [18]. The applied forecasting model for oil and gas power transformers obtained an acceptable solution with PSO [19], [20]. Hence, we improve the LSTM by integrating with PSO.

PSO star and ring topology, as well as a new particle representation that could improve the accuracy performance of crude oil forecasting, are embedded. The ring star topology and CHGSO_LSTM method, LSTM, and statistical time series techniques on benchmark crude oil price data compared the experimental results. Benchmark crude oil price data setting is the same as tested for CHGSO_LSTM by [14]. The rest of the paper follows the organization of the section as follows. Section II describes preliminaries on PSO and LSTM. The material and method for the proposed solution are in Section III. The computational results and discussion is

mentioned in Section IV and V respectively. Finally, Section VI concludes the paper.

## II. PRELIMINARIES

### A. Particle Swarm Optimization

In 1995, James Kennedy and Russell Eberhart introduced the PSO algorithm as a powerful population-based optimization technique [17], [18]. PSO has gained popularity among scientists and researchers due to its ease of implementation, high precision, and rapid convergence. The PSO algorithm is renowned for exploring and exploiting the search space effectively, rendering it suitable for various applications [17], [19]. PSO algorithm is a metaheuristic optimization technique that employs a population of particles to iteratively adjust their positions and velocities to find the best solution to a given problem. Before implementing PSO, particle representation must be designed carefully for the proper objective function [20]. The particle representation is an essential element of the PSO design for ensuring the algorithm's efficiency. Particle representation is a mechanism for encoding problem-solving solutions. Its ability to determine the properties of individual particles is used to map feature elements. By assigning an appropriate representation to each particle, PSO could facilitate efficient solutions [20], [21].

PSO can systematically investigate various regions within the search space while leveraging the search process to enhance and optimize a viable solution. The search strategies employed in the PSO algorithm are affected by the parameters, namely the acceleration constants ($C_1$ and $C_2$) and the inertia weight, as discussed by Shi and Eberhart [22]. Eq. (1) and Eq. (2) denote the velocity and position formulas adapted from the canonical PSO [18], [22].

$$V_{i(new)} = WV_i + C_1r_1(PBest_i - X_i) + C_1r_2(Pbest_i - X_i) + C_2r_2(Gbest_i - X_i) \quad (1)$$

$$X_{i(new)}) = X_i + V_{i(new)} \quad (2)$$

where,

$V_{i(new)}$ = new velocity of the $i^{th}$ particle

$V_i$ = current velocity of the $i^{th}$ particle

$X_i$ = current position of the $i^{th}$ particle

$X_{i(new)}$ = new position of the $i^{th}$ particle

$W$ = inertia weight

$C_1 \ and \ C_2$ = acceleration coefficient

$r_1$ and $r_2$ = random function in the range of [ 0, 1]

$PBest_i$ = position of the personal best of the $ith$ particle

$Gbest_i$ = position of the global best derived from all particles in the swarm

The cognitive component, represented as $C_1r_1(PBest_i - X_i)$ encompasses various factors such as the acceleration coefficient, a random function, and the difference between the personal best position, $PBest_i$, and the current position, $X_i$.

The difference between the previous position, $X_i$ and the personal best position $PBest_i$ can be observed. The social component, $C_2\, r_2\, (Gbest_i - X_i)$, incorporates the acceleration coefficient, a random function, and the disparity between the previous position, $X_i$, and the global best position, $Gbest_i$. The specific component signifies the historical performance of the particle, which is obtained from the combined version of all particles.

Another PSO strategy is topology. Topology is another PSO technique for exploiting and exploration. Topology controls how particles interact and exchange information, enabling them to jointly explore and seek the best outcome, such as the star, ring, and square. It establishes how information is shared and how interactions take place among particles to find the best solution.

### B. Long Term-Short Memory

LSTM can effectively capture long-range dependencies within sequential data [23]. This characteristic renders it highly appropriate for natural language processing, speech recognition, and time series analysis. Due to its feedback connections and capacity to acquire knowledge of long-time features from time series data, the LSTM network demonstrates significant efficacy in processing and predicting sequential data. The ability to capture extensive dependencies in sequential data is a considerable advantage of LSTM in deep learning [23].

LSTM is a deep learning technique demonstrating remarkable efficiency in capturing extensive dependencies within sequential data. LSTM is accomplished by employing memory cells alongside various gating mechanisms, including input, forget, and output gates. The architecture of the LSTM is depicted in Fig. 1. These mechanisms enable the extended short-term memory network to preserve and strategically discard information over duration, thereby facilitating its ability to accurately capture the interdependencies inherent in the dataset.

A neural network model's performance with dense and LSTM units depends on several variables, including the problem's complexity. Increasing the number of dense units in a neural network can enhance the model's ability to discover the patterns and relationships in the data when dealing with complex problems. Similarly, increasing the number of LSTM units would improve its ability to capture long-term dependencies and remember previous information over time. This is especially important in tasks involving sequential data, such as time series analysis.



Fig. 1. LSTM architecture.

Furthermore, dense units can transform input data into a higher-dimensional space, increasing the model's ability to separate and classify. However, it is crucial to exercise caution when selecting the appropriate number of dense and LSTM units to prevent the data's overfitting or underfitting. Furthermore, the structure and architecture of the neural network model can influence the effectiveness of dense and LSTM units. Considering the specific problem, it is recommended to experiment with different combinations of dense and LSTM units to find the optimal configuration that yields the best performance and generalization on the given task. The selection and number of dense and LSTM units during the construction of neural network models can impact the performance of the models [24] [25]. The configuration and selection of dense and LSTM units in a neural network model can significantly impact its performance.

### III. MATERIALS AND METHODS

This section elaborates on the description of the materials, data sources, and research methodologies used. The proposed approach captures the steps to see the performance of enhancement of LSTM with PSO models on crude oil forecasting. The approach includes data acquisition, pre-processing, construction of the proposed methods, and evaluation. We propose two variants of the LSTM+PSO model, including the LSTM+starPSO and LSTM+ringPSO models, and compare them with ARIMA, SARIMAX, and LSTM. Fig. 2 demonstrates the overview of the proposed methodology. In addition, we introduced a particle representation or solution mapping for the PSO. Detailed steps are elaborated in the following sub-sections.

### A. Data Acquisition

This study uses two different datasets: the WTI Crude Oil dataset [26] and the BRENT Crude Oil dataset [27]. Brent Crude refers to the assemblage of oil extracted from the North Sea's seabed, whereas WTI Crude denotes the amalgamation of oil obtained from land in the United States. WTI and BRENT are widely recognized benchmarks in the oil and gas industry. Specifically, the price of BRENT oil is commonly utilized as a reference point for the light oil market in Africa, Europe, and the Middle East. The datasets used for this study were obtained from the FRED website, a publicly accessible economic data repository owned by the Federal Reserve Bank of St. Louis. The website provides daily frequency data on WTI and BRENT crude oil prices from the early 1990s.

Nevertheless, the scope of this study is limited to the utilization of data solely from the period spanning from January 4, 2000, to April 15, 2021. The WTI dataset comprises 5409 objects, while the BRENT dataset contains 5438 objects. Both datasets share the same features: date, price, open, high, low, volume, and percentage change. A recent finding shows that the same dataset was used as a main part of a study by Altan and Karasu [28], which used two different forecasting methods, PSO+LSTM and CHGSO+LSTM. It was reported that the CHGSO+LSTM approach performed better than the LSTM method. The dataset's narrative was interestingly explored by using other LSTM variants.

Fig. 2.    Flow of methodology.

## B.  Data Cleaning

From the data behavior perspective, the two datasets, WTI and BRENT, exhibit no missing values. Therefore, no missing value procedure is imposed. However, the identification of outliers is required. The interquartile range (IQR) method can detect outliers by utilizing the interquartile range Data distribution should be determined in the interquartile range within Q1 and Q3 or between the 25th and 75th percentiles. An outlier is any data point that lies outside a predefined range, typically defined as below the 25th percentile and above the 75th percentile by about 1.5 times. This careful method of outlier discovery and eradication improves the overall data quality and ensures that the dataset is used for subsequent forecasting. Eq. (1) is the IQR formula [29].

$$IQR = Q3 - Q1 \qquad (3)$$

where, Q1 is the first quartile, and Q3 is the third quartile.

## C.  Proposed Method

This section explains an enhancement method of LSTM with PSO in forecasting crude oil. The initial part of the method construction is the identification of particle representation. A new particle representation is proposed to adhere to the LSTM architecture and its parameters. The aim is to find the best position of the particle that can give an optimal or near-optimal solution. It is represented by the particle's item, namely, lookback, LSTM unit, Dense Unit, and learning rate. The representation consists of discrete and continuous values shown in Fig. 3.

| Lookback | LSTM Unit | Dense Unit 1 | Dense Unit 2 | Learning Rate |
|---|---|---|---|---|
| [3 -10] | [64 – 256] | [10 -100] | [10 -100] | [0.01 - 0.001] |

Fig. 3.    Particle representation.

The LSTM-PSO algorithm incorporates two distinct topologies: the star and the ring. Consequently, two hybrid methodologies are proposed, specifically LSTM+starPSO and LSTM+ringPSO. Algorithm 1 outlines the procedural steps involved in the implementation of LSTM+starPSO. The LSTM+ringPSO follow similar steps, except Step 9. The algorithm commences by initializing the population of particles or swarm size. It is followed by initializing various parameters, including the number lookback, LSTM unit, dropout, dense unit, and learning rate. The Step 4 involves initializing the inertia weight, $w$, and acceleration constants $C_1$ and $C_2$ ). Steps 5 and 6 involve the initialization of the minimum value of velocity (Vmin), the maximum value of velocity (Vmax), the minimum position (Pmin), and the maximum value of position (Pmax). The subsequent step involves determining the setting for formulating the objective function and the iteration number, denoted as i. The ninth step involves the uploading of input data. Step 10 involves the implementation of the LSTM+starPSO algorithm, while Step 11 entails the computation of the Pbest and Gbest values for each particle. The updated characteristics of the particle are outlined in Step 14.

| **Algorithm 1: LSTM+starPSO** |
| --- |

| | |
| --- | --- |
| *1* | *Begin* |
| *2* | *Initialize the number of particles* |
| *3* | *Initiate number lookback, LSTM unit, dropout, dense unit, learning rate* |
| *4* | *Declare C1, C2 and W* |
| *5* | *Initialize Vmin, and Vmax* |
| *6* | *Initialize Pmax and Pmin.* |
| *7* | *Set the objective function based on RMSE* |
| *8* | *Set iteration number, i* |
| *9* | *Load datasets* |
| *10* | *Execute LSTM+starPSO* |
| *11* | *Calculate Pbest and Gbest value for each particle* |
| *12* | *Do* |
| *13* | *For each particle* |
| *14* | *Update features of the particle* |
| *15* | *Calculate the new velocity value, $V_{(new)}$* |
| *16* | *Calculate new position, $D_{(new)}$* |
| *17* | *Calculate Pbest $_{(new)}$* |
| *18* | *Calculate Gbest $_{(new)}$* |
| *20* | *While (stopping condition is reached)* |
| *21* | *End* |

The new velocity value for each particle is determined in Step 15 by applying Eq. (1). Eq. (2) is utilized to update the new position, denoted as P(new), in Step 16. Ultimately, Pbest(new) and Gbest(new) values are established by considering the fitness value assigned to the given problem. The iteration process commences at Step 12 and continues until Step 20, during which each particle's current velocity and position are updated. The iteration will continue until it meets the specified stopping condition.

### D. Performance Measure

This study has used two essential empirical measurements to evaluate and compare the effectiveness of the LSTM and PSO+LSTM models. The performance metrics, Mean Absolute Percentage Error (MAPE) and Root Mean Squared Error (RMSE) are the cornerstone indicators used to evaluate the precision and dependability of these models. RMSE is a well-known statistical metric that expresses the variance between the predicted values produced by the models and the actual observed values. A lower RMSE value signifies greater accuracy and precision, as the model's predictions closely match the observed data.

Second, by evaluating the relative error as a percentage of the actual values, MAPE provides an insightful perspective on the performance of the models. MAPE averages out the absolute percentage differences between the predicted and actual values. This metric is beneficial for Assessing how well the models can predict values roughly equivalent to the actual data points and approximately proportional to them. A lower MAPE indicates that the models make more accurate predictions with minor relative errors in applications.

## IV. COMPUTATIONAL RESULTS

### A. Parameter Setting

The proposed method encompasses two distinct parameter setting categories: Particle Swarm Optimization (PSO) and Long Short-Term Memory (LSTM). In the Particle Swarm Optimization (PSO) context, a population size for initializing particles is selected from a set of values, namely {10, 20, 30}. The value of the iteration variable, denoted as i, is adjusted to a value of 30. The importance of C1 and C2 remains consistently equal to 2. The lower and upper bounds for the inertia weight are 0.4 and 0.9, respectively. The particle representation set by PSO mapping determines the selection of random values for parameters such as lookback, LSTM unit, dense unit, and learning rate in the context of LSTM. The values are selected randomly during the execution of the program. The lookback parameter is randomly selected from 3 to 10, while the LSTM unit is chosen from 64 to 256. The density unit and learning rate values are specified within the range of [10, 100] and [0.01, 0.01], respectively.

### B. Computational Results using LSTM Based on the Number of Lookback

In assessing the impact of lookback on LSTM performance, we conducted experiments using varying lookback values, ranging from 3 to 10, as indicated in Table I. The objective is to determine an appropriate value for the lookback parameter and identify the optimal RMSE and MAPE values. As shown in Table I, the utilization of WTI in LSTM models yields diverse RMSE and MAPE values, indicating variations in performance. Two lookbacks, specifically 5 and 8, stand out due to their comparable RMSE and MAPE results. The root mean square error (RMSE) for Lookback 5 is calculated to be 2.6604, with a MAPE of 4.6256.

On the other hand, Lookback 8 has an RMSE of 2.7761 and a MAPE of 4.2628. Based on the analysis, it is observed that the MAPE of the lookback 8 model is significantly lower than that of the lookback 5 model, with a difference of -0.35. Additionally, the model with RMSE of lookback value equal to 8 is slightly higher than that of lookback equal to 5, with a difference of +0.11. Consequently, the model for lookback equal to 8 is deemed optimal for the WTI dataset. The results obtained on the BRENT dataset indicate a different outcome, highlighting the prominence of a particular lookback value of 7.

### C. Comparison Results of Different Methods

The forecasting results on the datasets WTI and BRENT are summarized in Table II and Table III. We incorporate LSTM with starPSO and ringPSO and tabulate the results from the recent finding by [28] and the conventional LSTM and two statistical time series methods, ARIMA and SARIMAX. The best results are highlighted in bold-face type. LSTM+starPSO provides better performance for the two datasets. In Table II, LSTM+starPSO with LSTM units of 212, dense unit of 77, and learning rate of 0.0083, lookback equals to 8 demonstrates the superior performance compared to other methods with RMSE of about 1.7512. We can see from Table III on the BRENT dataset a better forecasting

performance offered by LSTM+starPSO, where both RMSE and MAPE are minimized. The performance seems better than CHGSO_LSTM, which reported 1.7540 for WTI and 0.8453 of the RMSE for BRENT. In terms of the number of population, it shows that 20 is acceptable for both datasets. According to the results, although the data is univariate, the statistical models ARIMA and SARIMAX are ineffective compared to LSTM and its variants in forecasting future oil prices.

Every PSO requires an objective function or criterion that the PSO seeks to optimize. In this case, RMSE from the LSTM result is used as the objective function, which PSO tries to minimize the RMSE value. According to our findings, the ideal lookback range and learning rate are [6, 7] and 0.008,

respectively. LSTM+starPSO performs relatively similarly to the CHGSO-LSTM model [28] for the WTI and BRENT datasets in RMSE and MAPE. For the WTI dataset, +0.16% and -8.56% are obtained for RMSE and MAPE. A similar trend can be seen with the BRENT dataset, with +22.82% and +18.7% in RMSE and MAPE. On the other hand, LSTM+ringPSO performs slightly lower than LSTM+starPSO, where the result is -3.63% and -13.74% for RMSE and MPE for WTI and +20.7% and +21% for RMSE and MAPE for BRENT. PSO uses the position and velocity update method to find the best RMSE value. LSTM+starPSO outperforms CHGSO-LSTM with RMSE and MAPE with 11% and 5.7% reduction. The same goes for the BRENT dataset, where we see a 39% and 42.8% reduction in RMSE and MAPE, respectively.

TABLE I.    COMPUTATIONAL RESULTS USING CONVENTIONAL LSTM BASED ON THE NUMBER OF LOOKBACK

| Dataset | | WTI | | BRENT | |
|---|---|---|---|---|---|
| Split | Lookback | RMSE | MAPE | RMSE | MAPE |
| 60:40 | 3 | **2.1566** | **3.4105** | **2.4901** | **2.6015** |
| | 4 | 3.3722 | 4.7473 | 8.1081 | 7.1307 |
| | 5 | 4.2722 | 6.8263 | 7.2351 | 6.2737 |
| | 6 | 4.2524 | 6.3045 | 3.1759 | 3.0359 |
| | 7 | 10.7975 | 17.5782 | 5.2858 | 6.1048 |
| | 8 | 5.5147 | 8.2872 | 6.1838 | 7.7876 |
| | 9 | 5.2320 | 9.1245 | 4.1238 | 4.2867 |
| | 10 | 4.7698 | 5.5458 | 4.9142 | 4.9342 |
| 70:30 | 3 | 5.6011 | 10.4097 | 3.9737 | 5.3860 |
| | 4 | 4.1777 | 7.7009 | 3.3143 | 3.7884 |
| | 5 | 4.5288 | 7.8300 | **2.7562** | **3.2520** |
| | 6 | 5.6229 | 10.4402 | 3.1277 | 3.7722 |
| | 7 | **1.9410** | **3.4709** | 3.9283 | 4.7050 |
| | 8 | 3.2464 | 6.2048 | 5.4484 | 6.5811 |
| | 9 | 4.6569 | 8.9385 | 4.2039 | 5.3158 |
| | 10 | 4.3710 | 7.5679 | 2.7731 | 3.5558 |
| 85:15 | 3 | 2.6810 | 5.0949 | 3.0697 | 4.0615 |
| | 4 | 6.0544 | 10.7058 | 3.5370 | 4.9132 |
| | 5 | 2.6604 | 4.6256 | 3.2284 | 4.9966 |
| | 6 | 2.7652 | 5.0243 | 4.8779 | 7.9213 |
| | 7 | 3.2381 | 5.7339 | **2.1756** | **2.9873** |
| | 8 | **2.7761** | **4.2628** | 4.0058 | 6.2856 |
| | 9 | 3.3066 | 6.1230 | 2.7391 | 4.0462 |
| | 10 | 3.8408 | 7.4215 | 2.6276 | 3.9131 |
| 90:10 | 3 | 2.4417 | 4.8107 | 3.2680 | 3.7026 |
| | 4 | 5.1032 | 9.8931 | 2.5472 | 3.2090 |
| | 5 | 4.7821 | 8.8877 | <u>**1.8540**</u> | <u>**2.1971**</u> |
| | 6 | **2.1066** | **4.2554** | 4.6973 | 6.7682 |
| | 7 | 3.0119 | 6.0918 | 2.5806 | 3.2530 |
| | 8 | 3.4707 | 7.1888 | 3.5000 | 4.8287 |
| | 9 | 4.4182 | 8.6690 | 2.0140 | 2.5106 |
| | 10 | 4.0567 | 8.1125 | 2.3070 | 2.8610 |

TABLE II.     CRUDE OIL FORECASTING MODEL PERFORMANCE

| Performance Measure | CHGSO_LSTM [28] | LSTM+starPSO (Proposed method) | | | LSTM+ringPSO | | | LSTM | ARIMA | SARIMAX |
|---|---|---|---|---|---|---|---|---|---|---|
| Population | 20 | 20 | 10 | 30 | 20 | 10 | 30 | - | - | - |
| RMSE | 1.7540 | **1.7512** | 2.0601 | 1.9486 | 1.8199 | 2.0962 | 1.9889 | 1.9410 | 12.8125 | 15.0574 |
| MAPE | 2.7570 | **3.0150** | 3.5933 | 3.3466 | 3.1964 | 3.7775 | 3.4424 | 3.4709 | 24.8355 | 29.5014 |
| Learning Rate | 0.01 | **0.0083** | 0.0529 | 0.0076 | 0.0062 | 0.0470 | 0.0043 | 0.01 | 0.01 | 0.01 |
| LSTM units | 200 | **212** | 70 | 145 | 183 | 207 | 154 | 200 | 200 | 200 |
| Dense Units | 50 | **77** | 12 | 84 | 59 | 59 | 71 | 50 | 50 | 50 |
| Lookback | [3, 10] | **8** | 7 | 7 | 6 | 4 | 5 | 7 | - | - |

TABLE III.     BRENT CRUDE OIL FORECASTING MODEL PERFORMANCE

| Performance Measure | CHGSO_LSTM [28] | LSTM+star_PSO | | | LSTM+ringPSO | | | LSTM | ARIMA | SARIMAX |
|---|---|---|---|---|---|---|---|---|---|---|
| Swarm Size | 20 | 20 | 10 | 30 | 20 | 10 | 30 | - | - | - |
| RMSE | 1.8453 | **1.5024** | 2.0601 | 2.8145 | **1.5283** | 1.8231 | **1.7265** | 1.8540 | 14.9476 | 16.2943 |
| MAPE | 2.4525 | **2.0660** | **1.995**4 | 3.6217 | **2.0269** | **2.0359** | 2.2522 | 2.1971 | 25.1301 | 27.4191 |
| Learning Rate | 0.01 | 0.0066 | 0.0083 | 0.0011 | 0.0054 | 0.0051 | 0.0057 | 0.01 | 0.01 | 0.01 |
| LSTM Unit | 200 | 183 | 106 | 219 | 185 | 209 | 224 | 200 | - | - |
| Dense Unit | 50 | 51 | 60 | 38 | 65 | 27 | 14 | 50 | - | - |
| Lookback | [3, 10] | 10 | 3 | 5 | 7 | 7 | 7 | 5 | - | - |

## V. DISCUSSION

### A. Effect of LSTM and Dense Units

In LSTM, dense units can facilitate the transformation of the input data into a higher-dimensional space, increasing the model's ability to separate and classify. On the other hand, the number of dense and LSTM units should be chosen carefully to avoid overfitting or underfitting the data. Therefore, network architecture design, including the composition of dense and LSTM units, is significantly important [30]. It is advisable to experiment with different combinations of dense and LSTM units to find the optimal architecture that yields the best performance of forecasting accuracy. This paper explores using stochastic particle features to determine the most suitable number of dense and LSTM units and the incorporation of dense layers within an LSTM network. Interestingly, the LSTM+starPSO model, which used 212 LSTM units and 77 for dense units on WTI datasets, showed how sufficient dense and LSTM units reduce the overfitting problem and increase forecasting accuracy. In the context of the BRENT dataset, the model architecture consisted of 183 LSTM units and 51 dense units.

### B. Effect of Lookback

Lookback in time series forecasting establishes how much historical data the LSTM, LSTM+starPSO, and LSTM+ringPSO models should consider when making predictions for the following time step. Depending on several variables, including the data patterns, the effect of lookback on LSTM can be significant. Determining the most effective lookback period is contingent upon the unique attributes of the time series data [31]. More than eight lookback numbers appear required for accurate prediction in the best performance models, which offer little but longer-term dependencies. With such a small number of lookbacks, more is needed.

### C. Effect of Learning Rate

The role of the learning rate in the LSTM forecasting model is to assist in effective converging and achieving good performance. We represent the particle with the range of learning rate of [0.01, 0.001]. It is randomly chosen within these ranges during the LSTM+starPSO and LSTM+ringPSO execution. The small learning rate value is randomly chosen at about 0.006 to 0.008 for both datasets using LSTM+starPSO, meanwhile about 0.004 to 0.005 when using LSTM+ringPSO. However, the use of a small learning rate value has a significant effect on the forecasting results. It is evident that the choice of a small learning rate in LSTM models for time series forecasting aims to achieve a better convergence and generalization in the exploitation and exploration of the search space [32].

### D. Effect of PSO Topology in LSTM

Ring topology is a local-based focal point of particles. It attracts particles to the best particle in its corresponding neighborhood. In our experiment, three swarm sizes are used: 10, 20, 30. For instance, each particle has 29 neighborhoods when using a swarm size equal to 30. However, due to the various particle positions in the search space, the nearest particle of each particle considers the local neighborhood involved in the local search. Each particle's local surroundings consist of a fixed number of other particles. It differs from star topology, where all particles within the swarm share information with and are influenced by the particle with the highest performance [33]. Star topology promotes global exploration by encouraging particles to move towards the

optimal solution discovered by any swarm member. LSTM with star topology has demonstrated that each particle in the searching space is attracted to the best particle of the swarm. It obtained the best forecasting accuracy performance for WTI and BRENT datasets. The global searching by star topology [34] achieves the objective function of minimizing the RMSE value.

## VI. CONCLUSION

In this study, an enhancement of LSTM incorporated with PSO addresses the challenges of forecasting the daily time series crude oil price data. The proposed method comprises two main steps. At the PSO, particle mapping is designed together with topology to achieve a dynamic LSTM architecture and improve PSO searching capabilities for exploration and exploitation. With this method, more accurate forecasting is obtained. Experimental findings show that compared with the recent CHGSO_LSTM, the suggested LSTM+starPSO offers the most performing methods. It is a better outcome compared to LSTM+ringPSO and conventional methods. However, more experimental work could be conducted by embedding ensembles and executing feature engineering strategy and hyperparameter tuning.

## ACKNOWLEDGMENT

## REFERENCES

[1] Z. Xu, M. Mohsin, K. Ullah, and X. Ma, "Using econometric and machine learning models to forecast crude oil prices: Insights from economic history," Resources Policy, vol. 83, p. 103614, 2023, doi: 10.1016/J.RESOURPOL.2023.103614.

[2] S. Lipsa and R. K. Dash, "GASVR- A Model to Predict and Analyze Crude Oil Price," in 2022 2nd Asian Conference on Innovation in Technology, ASIANCON 2022, 2022. doi: 10.1109/ASIANCON55314.2022.9908764.

[3] H. Son and Y. Jang, "Partial convolutional LSTM for spatiotemporal prediction of incomplete data," IEEE Access, vol. 8, pp. 164762–164774, 2020, doi: 10.1109/ACCESS.2020.3022774.

[4] P. Sohrabi, H. Dehghani, and R. Rafie, "Forecasting of WTI crude oil using combined ANN-Whale optimization algorithm," Energy Sources, Part B: Economics, Planning and Policy, vol. 17, no. 1, 2022, doi: 10.1080/15567249.2022.2083728.

[5] M. E. A. Ben Seghier, D. Höche, and M. Zheludkevich, "Prediction of the internal corrosion rate for oil and gas pipeline: Implementation of ensemble learning techniques," J Nat Gas Sci Eng, vol. 99, p. 104425, 2022, doi: https://doi.org/10.1016/j.jngse.2022.104425.

[6] M. Farsi et al., "Predicting Formation Pore-Pressure from Well-Log Data with Hybrid Machine-Learning Optimization Algorithms," Natural Resources Research, vol. 30, no. 5, pp. 3455–3481, 2021, doi: 10.1007/s11053-021-09852-2.

[7] K. Tissaoui, T. Zaghdoudi, A. Hakimi, and M. Nsaibi, "Do Gas Price and Uncertainty Indices Forecast Crude Oil Prices? Fresh Evidence Through XGBoost Modeling," Comput Econ, vol. 62, no. 2, pp. 663–687, 2023, doi: 10.1007/s10614-022-10305-y.

[8] X.-L. Pan, L.-J. Zhang, G. He, and P. Zhang, "Prediction model for residual strength of warship seawater pipelines with internal corrosion and test verification," Chuan Bo Li Xue/Journal of Ship Mechanics, vol. 25, no. 2, pp. 202–209, 2021, doi: 10.3969/j.issn.1007-7294.2021.02.008.

[9] R. Yang, X. Liu, R. Yu, Z. Hu, and X. Duan, "Long short-term memory suggests a model for predicting shale gas production," Appl Energy, vol. 322, 2022, doi: 10.1016/j.apenergy.2022.119415.

[10] V. Shahbazbegian, H. Hosseininesaz, M. Shafie-Khah, and M. Elmusrati, "Forecasting Crude Oil Prices using a Hybrid Model Combining Long Short-Term Memory Neural Networks and Markov Switching Model," in 2023 International Conference on Future Energy Solutions, FES 2023, 2023. doi: 10.1109/FES57669.2023.10182444.

[11] S. Song et al., "An intelligent data-driven model for virtual flow meters in oil and gas development," Chemical Engineering Research and Design, vol. 186, pp. 398–406, 2022, doi: 10.1016/j.cherd.2022.08.016.

[12] A. S. Dyer et al., "Applied machine learning model comparison: Predicting offshore platform integrity with gradient boosting algorithms and neural networks," Marine Structures, vol. 83, 2022, doi: 10.1016/j.marstruc.2021.103152.

[13] H. He, M. Sun, X. Li, and I. A. Mensah, "A novel crude oil price trend prediction method: Machine learning classification algorithm based on multi-modal data features," Energy, vol. 244, p. 122706, 2022, doi: 10.1016/J.ENERGY.2021.122706.

[14] A. Altan and S. Karasu, "Crude oil time series prediction model based on LSTM network with chaotic Henry gas solubility optimization," Energy, vol. 242, p. 122964, 2022, doi: 10.1016/J.ENERGY.2021.122964.

[15] H. Guliyev and E. Mustafayev, "Predicting the changes in the WTI crude oil price dynamics using machine learning models," Resources Policy, vol. 77, 2022, doi: 10.1016/j.resourpol.2022.102664.

[16] H. He, M. Sun, X. Li, and I. A. Mensah, "A novel crude oil price trend prediction method: Machine learning classification algorithm based on multi-modal data features," Energy, vol. 244, 2022, doi: 10.1016/j.energy.2021.122706.

[17] M. Yusoff, J. Ariffin, and A. Mohamed, "DPSO based on a min-max approach and clamping strategy for the evacuation vehicle assignment problem," Neurocomputing, vol. 148, pp. 30–38, 2015, doi: https://doi.org/10.1016/j.neucom.2012.12.083.

[18] X. Hu, Y. Shi, and R. Eberhart, "Recent advances in particle swarm," in Proceedings of the 2004 Congress on Evolutionary Computation, CEC2004, 2004, pp. 90–97. [Online]. Available: https://www.scopus.com/inward/record.uri?eid=2-s2.0-4344682502&partnerID=40&md5=60528ec5f19a0f07655319c9c06b8e28.

[19] M. Riaz, A. Hanif, S. J. Hussain, M. I. Memon, M. U. Ali, and A. Zafar, "An optimization-based strategy for solving optimal power flow problems in a power system integrated with stochastic solar and wind power energy," Applied Sciences (Switzerland), vol. 11, no. 15, 2021, doi: 10.3390/app11156883.

[20] Y. Shi and R. Eberhart, "Monitoring of Particle Swarm Optimization," Front Comput Sci China, vol. 3, no. 1, pp. 31–37, 2009, doi: 10.1007/s11704-009-0008-4.

[21] M. Yusoff, A. N. M. Basir, N. A. Kadir, and S. A. Bahari, "Evaluation of Particle Swarm Optimization for strength determination of tropical wood polymer composite," IAES International Journal of Artificial Intelligence, vol. 9, no. 2, pp. 364–370, 2020, doi: 10.11591/ijai.v9.i2.pp364-370.

[22] Y. Shi and R. C. Eberhart, "Population diversity of particle swarms," in 2008 IEEE Congress on Evolutionary Computation, CEC 2008, 2008, pp. 1063–1067. doi: 10.1109/CEC.2008.4630928.

[23] X. Lei, "Stock Market Forecasting Method Based on LSTM Neural Network," in 2023 IEEE 3rd International Conference on Power, Electronics and Computer Applications, ICPECA 2023, 2023, pp. 1534–1537. doi: 10.1109/ICPECA56706.2023.10076100.

[24] K. Gajamannage, Y. Park, and D. I. Jayathilake, "Real-time forecasting of time series in financial markets using sequentially trained dual-LSTMs," Expert Syst Appl, vol. 223, 2023, doi: 10.1016/j.eswa.2023.119879.

[25] M. S. Sawah, S. A. Taie, M. H. Ibrahim, and S. A. Hussein, "An accurate traffic flow prediction using long-short term memory and gated recurrent unit networks," Bulletin of Electrical Engineering and Informatics, vol. 12, no. 3, pp. 1806–1816, 2023, doi: 10.11591/eei.v12i3.5080.

[26] "Crude Oil Price | WTI Price Chart - Investing.com." 2023. [Online]. Available: https://www.investing.com/commodities/crude-oil.

[27] "Brent Crude Oil Price - Investing.com." 2023. [Online]. Available: https://www.investing.com/commodities/brent-oil.

[28] A. Altan and S. Karasu, "Crude oil time series prediction model based on LSTM network with chaotic Henry gas solubility optimization," Energy, vol. 242, p. 122964, 2022, doi: 10.1016/J.ENERGY.2021.122964.

[29] K. Fujioka and M. Deng, "Nonlinear temperature control of heat exchanger process using particle filter with interquartile range particle weights," International Conference on Advanced Mechatronic Systems, ICAMechS, vol. 2021-December, pp. 130–134, 2021, doi: 10.1109/ICAMECHS54019.2021.9661547.

[30] A. A. Ismail and M. Yusoff, "An Efficient Hybrid LSTM-CNN and CNN-LSTM with GloVe for Text Multi-class Sentiment Classification in Gender Violence," International Journal of Advanced Computer Science and Applications, vol. 13, no. 9, 2022, doi: 10.14569/IJACSA.2022.0130999.

[31] H. Jeon, J. Ryu, K. M. Kim, and J. An, "The Development of a Low-Cost Particulate Matter 2.5 Sensor Calibration Model in Daycare Centers Using Long Short-Term Memory Algorithms," Atmosphere (Basel), vol. 14, no. 8, 2023, doi: 10.3390/atmos14081228.

[32] C. Yu, X. Qi, H. Ma, X. He, C. Wang, and Y. Zhao, "LLR: Learning learning rates by LSTM for training neural networks," Neurocomputing, vol. 394, pp. 41–50, 2020, doi: https://doi.org/10.1016/j.neucom.2020.01.106.

[33] N. Lynn, M. Z. Ali, and P. N. Suganthan, "Population topologies for particle swarm optimization and differential evolution," Swarm Evol Comput, vol. 39, pp. 24–35, 2018, doi: https://doi.org/10.1016/j.swevo.2017.11.002.

[34] V. Mann, A. Sivaram, L. Das, and V. Venkatasubramanian, "Robust and efficient swarm communication topologies for hostile environments," Swarm Evol Comput, vol. 62, p. 100848, 2021, doi: https://doi.org/10.1016/j.swevo.2021.100848.

# Explore Innovative Depth Vision Models with Domain Adaptation

Wenchao Xu[1], Yangxu Wang[2]

School of Electrical and Computer Engineering, Nanfang College Guangzhou, Conghua 510970, China[1]

Department of Network Technology, Software Engineering Institute of Guangzhou, Conghua 510990, China[2]

*Abstract*—In recent years, deep learning has garnered widespread attention in graph-structured data. Nevertheless, due to the high cost of collecting labeled graph data, domain adaptation becomes particularly crucial in supervised graph learning tasks. The performance of existing methods may degrade when there are disparities between training and testing data, especially in challenging scenarios such as remote sensing image analysis. In this study, an approach to achieving high-quality domain adaptation without explicit adaptation was explored. The proposed Efficient Lightweight Aggregation Network (ELANet) model addresses domain adaptation challenges in graph-structured data by employing an efficient lightweight architecture and regularization techniques. Through experiments on real datasets, ELANet demonstrated robust domain adaptability and generality, performing exceptionally well in cross-domain settings of remote sensing images. Furthermore, the research indicates that regularization techniques play a crucial role in mitigating the model's sensitivity to domain differences, especially when incorporating a module that adjusts feature weights in response to redefined features. Moreover, the study finds that under the same training and validation set configurations, the model achieves better training outcomes with appropriate data transformation strategies. The achievements of this research extend not only to the agricultural domain but also show promising results in various object detection scenarios, contributing to the advancement of domain adaptation research.

*Keywords—Deep learning; neural network; domain adaptation; lightweight; regularization techniques*

## I. INTRODUCTION

Remote sensing technology, as a crucial tool for observing the Earth and its ecosystems, has played an indispensable role in multiple domains [1]. However, due to various factors such as capture devices, time, and location influencing the acquisition process of remote sensing images, there exist differences in remote sensing data across different domains. Consequently, models trained in one domain (source domain) often exhibit decreased performance when applied to another domain (target domain). This challenge is commonly referred to as distributional difference [2].

Indeed, distributional difference has long been a persistent issue in machine learning. A series of studies have demonstrated that as the mismatch between distributions increases, performance noticeably declines [3]. A widely adopted approach to address this issue is Domain Adaptation (DA). Previous research has shown that domain adaptation markedly impacts the accuracy and reliability of processing remote sensing images. To tackle this problem, researchers have explored various methods, such as reannotating a portion of target domain data for model fine-tuning [2], making the data distributions of the source and target domains more similar through feature selection or transformation [4], utilizing adversarial training for domain alignment [5], and employing strategies like self-supervised learning [6] and meta-learning [7]. These approaches have, to some extent, alleviated the problem of distributional differences.

With the rapid development of the third wave of artificial intelligence—deep learning, deep convolutional neural networks (CNNs) are significantly pushing the performance boundaries of computer vision at an incredible pace [8]. The latest advances in Unsupervised Domain Adaptation (UDA) in image processing have been attempted and progressed in various fields. Goel et al. [9] achieved unsupervised domain adaptation by guiding transfer learning and employing the Jensen-Shannon (JS) divergence method. In the remote sensing domain, Elshamli et al. [10] introduced an innovative approach to domain adaptation, incorporating denoising autoencoders and domain adversarial neural networks, especially in the classification of hyperspectral and multispectral images. In the agricultural domain, Zhang et al. [11] narrowed the gap between the source and target domains for agricultural land extraction using Generative Adversarial Networks (GANs). Similarly, Valerio et al. [12] accomplished unsupervised leaf counting, while Marino et al. [13] achieved potato defect classification. In plant disease recognition, Fuentes et al. [14] proposed open-set adaptation and cross-domain adaptation methods to enhance tomato disease recognition using unlabeled data. Additionally, Wu et al. [15] achieved cross-domain recognition of wild plant diseases. In robotics, Magistri et al. [16] introduced Unsupervised Domain Adaptation (UDA) techniques for semantic segmentation, enhancing the adaptability of agricultural robots to better perceive and understand different environments. These collective efforts address challenges associated with domain transfer and variability, contributing to the robustness and adaptability of image processing models for various application domains.

Despite these advancements, it is noteworthy that the visual representations learned by deep CNNs exhibit considerable domain invariance. There is evidence that combining existing CNN representations with a linear classifier can achieve relatively high accuracy [17]. In earlier research, Lu et al. [18] posed a challenging question, namely achieving domain adaptation without explicit adaptation of data distribution. This

provided an opportunity to reconsider the problem of cross-domain generalization from a new perspective.



Fig. 1. Two typical scenarios. The target domain is shown in gray and the source domain is shown in red. Different tags indicate different categories.

As Fig. 1 illustrates two typical data distributions, an interesting phenomenon can be observed: samples from different domains but of the same class are close enough, while samples from different classes have a sufficiently large gap. In scenario (a), the distribution difference between the source and target domains is small, but due to confusion among different classes, the final recognition performance is not ideal. In contrast, in scenario (b), despite a marked difference between different domains, samples from different classes are still linearly separable. This suggests that the magnitude of differences between the source and target domains is not a good indicator of the final recognition accuracy. Is the CNN representation powerful enough to eliminate the need for domain adaptation? This study aim to explore a visual model that achieves domain adaptation without the need for explicit adaptation, designed a novel deep convolutional neural network, Efficient Lightweight Aggregation Network (ELANet), which innovatively employs an ELA module for optimizing feature decoding. Experimental results demonstrate the significant optimization of ELANet, showcasing domain adaptability due to the powerful representational capabilities of current CNNs and some carefully designed features. The validation process is based on three datasets: the Global Wheat Head Detection 2021 (GWHD) dataset focusing on wheat [19], and two remote sensing datasets, namely Remote Sensing Object Detection (RSOD) [20] and University of Chinese Academy of Sciences - Aerial Object Detection (UCAS-AOD) [21]. For the remote sensing datasets, the "airplane" category was selected, with RSOD serving as the source domain and UCAS-AOD as the target domain. Extensive experimental results demonstrate the effectiveness of the ELANet method, and some interesting findings are reported.

In summary, the contributions of the research can be summarized as follows:

- ELANet: A visual model with domain adaptation, reporting state-of-the-art performance in cross-domain settings for agricultural and remote sensing scenarios, demonstrating sufficient generality.

- Validation of the effectiveness of regularization techniques: The study proves that utilizing regularization techniques to mitigate domain variance is

effective. In particular, incorporating a module that dynamically adjusts feature weights in response to domain redefinition is a more intelligent approach.

- Effective Training Set and Validation Set Configuration: Training under the same configuration for training and validation sets is more effective, provided the existence of data transformation strategies.

## II. METHODOLOGY

### A. ELANet Model Design

The ELANet model is designed based on two components: Encoder and Decoder. The Encoder extracts and downsamples features from input images through a series of convolutional and channel transformation layers, forming feature maps at different resolutions. The Decoder is responsible for the feature extraction task and includes multiple branches. Each branch processes features at different scales through convolutional and channel fusion operations. Multi-scale feature fusion enhances detection performance. Finally, an Adaptive Scale Fusion (ASF) layer is employed to adaptively fuse features, reducing the need for deep downsampling to obtain high-level semantic information [22], as depicted in Fig. 2. The following sections will introduce the global architecture of the ELANet model and its optimizations.



Fig. 2. The architecture of ELANet.



Fig. 3. Details of module design.

### B. Encoder

The role of the encoder is to map the input RGB images into feature maps. Specifically, includes five downsampling operations performed by convolutional layers with a stride of 2 and a 3×3 kernel size. In the middle, a C2f module [23] is inserted for feature extraction, generating feature maps at different stages. The sequence of operations can be described as follows: Conv3-Conv3-C2f-Conv3-C2f-Conv3-C2f-Conv3-C2f (In Convk, 'k' represents the size of the convolution kernel). These operations are connected sequentially, with the

output of the previous layer serving as the input for the next layer. The initial channel number is 3, corresponding to the RGB image channels, and it increases gradually. It's noteworthy that in the final stage of the Encoder, a Multi-Head Self-Attention module called Attention-based Intra-scale Feature Interaction (AIFI) [24] is employed to handle the highest-level features of the backbone network, as depicted in Fig. 3. The mathematical processes are defined by Eq. (1) and Eq. (2):

$$Q = K = V = Flatten(S_{Last}) \qquad (1)$$

$$F_{Last} = Reshape(Attn(Q, K, V)) \qquad (2)$$

where, $S_{Last}$ represents the last layer feature map output by the Encoder. Initially, the two-dimensional feature $S_{Last}$ is flattened into a vector, which is then processed by the AIFI module. Subsequently, the output is reshaped back into two dimensions, denoted as $F_{Last}$, facilitating its transmission to the Decoder for feature analysis.

*C. Decoder*

The role of the Decoder is to combine and utilize features from the Encoder and decode predictive information. In visual models for processing remote sensing images, the utilization of gradient information is crucial. This process is generally achieved through continuous upsampling to facilitate the fusion of semantic information. Despite being a common approach, there is inevitably a problem of feature information loss or degradation, impacting the fusion effectiveness across non-adjacent levels. To address this issue, this study introduces the ELA module, whose structure is depicted in Fig. 3. Specifically, the ELA module employs two consecutive convolutional operations to capture richer feature representations. Subsequently, by introducing additional gradient flow branches in parallel, incorporates a broader context to enhance abstraction capabilities for targets. The features from the first six layers are concatenated to simultaneously utilize multi-scale information, strengthening the detection capabilities for targets of different sizes. Finally, a 1×1 convolutional layer is applied to reduce the dimensionality of the matrix, thereby alleviating the computational load of the model. With these connection operations, the model gains richer gradient information, achieving higher accuracy and more reasonable latency. Notably, in the ELANet model, a further adaptation is made using Adaptive Spatial Fusion (ASF) proposed by Yang et al. [25], which supports direct interaction between non-adjacent levels for adaptive spatial feature fusion. Through two consecutive convolutional operations, the ELA module adapts well to objects of different scales. The first convolutional operation captures features at a smaller scale, while the second convolutional operation integrates these features within a larger receptive field, making the model more flexible and accurate in detecting objects of different sizes. The design of the Decoder section integrates features from three different levels of the Encoder. ASF allocates different spatial weights to enhance the importance of key levels and reduce the impact of conflicting information from different levels.

Representation learning in convolutional neural networks faces the challenge of strong correlations between adjacent pixels, implying the potential provision of redundant information. To address this issue, there are some carefully designed strategies within each output feature of ASF, utilizing regularization techniques to alleviate domain variance. Specifically, a dropout strategy with a probability of 0.1 and Shuffle Attention (SA) attention strategy [26] during the upsampling and downsampling processes are employed. Since the zeroing of elements after dropout is random, connecting an SA attention strategy for responsive feature weight adjustment is deemed necessary. This necessity will be validated and analyzed in the experiments below. Finally, ELANet merges the decoded feature maps from different stages into the original feature image. Through pixel-level prediction, the regression branch predicts the distance from each anchor point to the four edges of the target bounding box, determining the target's position.

*D. Loss Function*

In the implementation of the ELANet model, use three loss functions to guide the regression of bounding boxes. The Classification Loss is used to measure the difference between the predicted class and the true class. This process is guided by the cross-entropy loss function, a common binary classification loss function, defined as follows:

$$L_c = -\frac{1}{n}\sum_{i=1}^{n}[y_i \log p_i + (1 - y_i)\log(1 - p_i)] \qquad (3)$$

In Eq. (3), let the ground truth be denoted as y, the predicted result as y, and n represents the batch size.

The Regression Loss is composed of Complete Intersection over Union (CIoU) and Distribution Focal Loss (DFL). CIoU is used to guide the model in learning the matching degree of bounding boxes. Specifically, assuming the predicted box and the ground truth box are denoted as $b_p$ and $b_g$ respectively, it is described as follows:

$$L_{CIoU} = IoU - \frac{\rho^2(b_p, b_g)}{c^2} - \alpha v \qquad (4)$$

In Eq. (4), IoU represents the Intersection over Union, and $\rho^2(b_p, b_g)$ calculates the Euclidean distance between the center points of the two rectangular boxes. Here, $c$ is the length of the diagonal of the minimum bounding rectangle of the predicted and ground truth boxes, $v$ is used to measure the similarity of aspect ratios, and $\alpha$ is the impact factor of $v$. Next, $dfl$ optimizes the position of the bounding boxes through smooth L1 loss. For each positive sample $i$, it is defined as:

$$L_{dfl} = \frac{1}{N_{pos}}\sum_{i=1}^{N_{pos}} Sml(pd_i, gt_i) \times w_i \qquad (5)$$

In Eq. (5), $N_{pos}$ represents the number of positive samples, $Sml(pd_i, gt_i)$ represents the smooth L1 loss for the i[th] positive sample, and $w_i$ is the weight associated with the i[th] sample. Combining Eq. (4) and Eq. (5), the regression loss Lr can be obtained as $L_r = L_{CIoU} + L_{dfl}$.

The final loss for ELANet is defined as the weighted sum of the classification loss and regression loss, i.e., $L_{ELANet} = \alpha L_c + \beta L_r$.

## III. EXPERIMENT

### A. Experimental Details

*1) Data preprocessing:* The RSOD and UCAS-AOD datasets consist of 446 and 1000 airplane images, respectively, designated as the source domain and target domain. For convenience, this configuration is named RSOD. Beyond that, the GWHD dataset was intentionally set up for cross-domain settings, with 165 wheat spike images in the source domain and 71 in the target domain. In Fig. 4, present the size information for each instance in these two datasets. The RSOD dataset exhibits significant size differences between source and target domain data, while the primary difference in the GWHD dataset originates from variations in the external and internal environments.

*2) Training details:* The experiments are implemented in PyTorch [27] and accelerated using an NVIDIA RTX 3090 GPU. To improve computational efficiency, the longest side of input images is scaled to 608 pixels, and the other side is scaled proportionally, which also suits the resolution requirements when deploying on low-end edge devices. The Adaptive Moment Estimation (Adam) algorithm [28] was employed as the optimizer with a momentum factor set to 0.937, and the initial learning rate was set to 0.01. Considering convergence speed, perform 150 epochs of optimization on the RSOD dataset and 300 epochs on the GWHD dataset. To ensure the robustness of model training, strategies such as color distortion, random scale transformations, and mosaic data augmentation are employed.



Fig. 4. Instance size information for each dataset.

### B. Evaluation Indicators

In the process of establishing a detection model, it is essential to consider both precision and recall. Therefore, this study adopts metrics such as Precision, Recall, mAP@0.5, and mAP@0.5-0.95 to evaluate the model's performance and assess the detection results. The calculation methods for Precision and Recall are as Eq. (6) and Eq. (7):

$$P = \frac{TP}{TP+FP} \qquad (6)$$

$$R = \frac{TP}{TP+FN} \qquad (7)$$

where, P represents precision, and R represents recall. True Positive (TP) is the number of positive samples correctly classified, True Negative (TN) is the number of negative samples correctly classified, False Positive (FP) is the number of negative samples incorrectly classified as positive, and False Negative (FN) is the number of positive samples incorrectly classified as negative.

mAP represents the comprehensive performance at different Intersection over Union (IoU) thresholds, including mAP@0.5 and mAP@0.5-0.95. Here, mAP@0.5 denotes the average mAP when the IoU threshold is 0.5, with a higher value indicating higher detection precision for that category. mAP@0.5-0.95 represents the average mAP at different IoU thresholds (ranging from 0.5 to 0.95 with a step size of 0.05), placing stricter demands on the model's performance. The calculation of mAP is given by Eq. (8):

$$mAP = \frac{1}{n}\sum_1^n P(R)d(R) \qquad (8)$$

where, n is the number of categories. In this experiment, each dataset has only one label, so n=1.

Furthermore, three metrics will be employed in the counting evaluation to assess the consistency between predicted and ground truth values, including Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and Coefficient of Determination (R²). Specifically, they are defined by Eq. (9) to Eq. (11):

$$MAE = \frac{1}{n}\sum_{i=1}^n|\widehat{y_i} - y_i| \qquad (9)$$

$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^n(\widehat{y_i} - y_i)^2} \qquad (10)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n(\widehat{y_i}-y_i)^2}{\sum_{i=1}^n(\overline{y_i}-y_i)^2} \qquad (11)$$

In these formulas, the numerator represents the sum of squared differences between the actual values and predicted values, and the denominator represents the sum of squared differences between the actual values and the mean. The result of $R^2$ falls within the range of [0, 1], indicating the proportion of the squared differences of predicted values to the squared differences of actual values near the mean. This metric can be understood as a measure of how well the model's predictions fit the actual values, with 1 indicating a perfect fit and 0 indicating no linear relationship between actual counts and predicted values.

### C. Comparision with other Methods

To validate the effectiveness of ELANet, comparisons with two state-of-the-art methods: the two-stage model Faster R-CNN [29] and the one-stage model YOLOv7-tiny [30], both of which are widely used for visual tasks in images. Table I and Table II present the quantitative results on the two datasets, while Fig. 5 showcases some prediction examples from ELANet.

One notable observation is that, even though both datasets explicitly implement cross-domain settings, the performance on RSOD significantly surpasses that on GWHD. In reality, domain differences in the agricultural domain are extremely complex. The chaotic background makes the visual patterns of plants diverse and misleading, and the changes in plants themselves are also very pronounced. As shown in Fig. 1(a),

despite the small distribution difference between the source and target domains, the recognition performance is not ideal due to the confusion of different category samples. As demonstrated in Fig. 5(a), (b), the morphological differences in wheat spikes from different regions and varieties are substantial.



Fig. 5. Part of ELANet's prediction results. (a, b) are from the GWHD dataset, and (c, d) are from the RSOD dataset.

The comparison of object detection methods on the GWHD and RSOD datasets reveals distinct performance characteristics. According to the quantitative results shown in Table I for models on the GWHD dataset, ELANet achieves a precision of 88.7% without the need for adaptation to distribution differences on low-resolution input images. YOLOv7-tiny, while effective in terms of parameters, exhibits moderate precision and recall, with a lower mAP@0.5-0.95 score. Faster R-CNN, with a higher parameter count, shows performance comparable to YOLOv7-tiny but lacks the precision and recall achieved by ELANet, possibly due to information loss resulting from the simplification of their network structures. ELANet achieves optimal performance at the same input size of 608×608, substantially improving precision and average precision.

TABLE I. QUANTITATIVE RESULTS OF GWHD DATASET

| Method | P | R | mAP@0.5 | mAP@0.5-0.95 | Params |
|---|---|---|---|---|---|
| YOLOv7-tiny | 0.558 | 0.379 | 0.387 | 0.141 | 11.55M |
| Faster R-CNN | 0.446 | 0.391 | 0.343 | 0.127 | 42.20M |
| ELANet | 0.882 | 0.816 | 0.887 | 0.501 | 3.59M |

Furthermore, as shown in Table II, the experimental results on the RSOD dataset also demonstrate a similar trend. However, the evaluation gap between the models is smaller in the RSOD dataset experiments. Additionally, it is worth noting that ELANet exhibits a smaller model parameter size, only 3.59M, compared to YOLOv7-tiny and Faster R-CNN with 11.55M and 42.2M, respectively. ELANet achieves significant improvements in both performance and parameter efficiency, indicating that it maintains high performance while being more parameter-efficient, making it well-suited for lightweight tasks without sacrificing efficiency.

TABLE II. QUANTITATIVE RESULTS OF RSOD DATASET

| Method | P | R | mAP@0.5 | mAP@0.5-0.95 | Params |
|---|---|---|---|---|---|
| YOLOv7-tiny | 0.555 | 0.821 | 0.802 | 0.310 | 11.55M |
| Faster R-CNN | 0.832 | 0.751 | 0.753 | 0.278 | 42.20M |
| ELANet | 0.894 | 0.845 | 0.861 | 0.337 | 3.59M |

*D. Analysis of Counting Influencing Factors*

Due to the dense distribution of targets in both the GWHD and RSOD datasets, evaluating counting performance in this context is meaningful. Particularly in the case of the GWHD wheat head dataset, the model is prone to various natural factors during target detection, such as the influence of lighting, rainy weather, thick fog, etc. What's more, errors may arise from the varying shapes, sizes, and densities of different wheat head varieties. In this experiment, a linear regression plot was employed, along with counting metrics introduced in Section III (B), including Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and the Coefficient of Determination ($R^2$), to assess counting performance. The diagonal line on the coordinate system represents the ideal state where the model inference results perfectly match the manually counted ground truth. The linear regression plot serves as an effective tool to visually analyze the relationship between model predictions and actual counts, allowing researchers to gain a deeper understanding of the factors influencing counting performance. Additionally, it highlights images with the highest errors, as shown in Fig. 6.

It is evident that the maximum errors are primarily concentrated under varying lighting conditions, where inconsistencies in illumination affect wheat heads in bright and shadowed areas differently. The presence of cluttered foliage further challenges the target detection model. Through experimentation, it was discovered that even human experts find it challenging to discern in such conditions. Despite this, the ELANet model demonstrates good robustness, indicating that the optimized strategies of data augmentation play a positive role in enhancing the adaptability of the model for target detection performance. This also points towards future research directions, specifically addressing how to optimize models for strong interference environments to achieve stability, reliability, and adaptability in complex conditions.



Fig. 6. Linear regression plot and maximum error plot of GWHD dataset count results.

*E. Ablation Study*

In the RSOD dataset, ELANet's performance is satisfactory. Building upon this, the focus shifted to conducting ablation experiments using Dropout strategy and SA strategy on the GWHD dataset. The experimental results are presented

in Table III and compared with two domain adaptation methods. In the table, "Dp" represents the use of the dropout strategy, and "At" represents the application of the SA attention strategy. It can be observed that without employing either strategy, ELANet achieves an accuracy of 0.87. However, using either the dropout strategy or attention strategy alone leads to a performance decrease in the detection task, with reductions of 1.72% and 0.23%, respectively. On the other hand, employing the SA attention strategy after using the dropout strategy proves to be highly effective. The random zeroing of elements after dropout reduces the risk of overfitting, preventing excessive co-adaptation of neurons and maintaining feature diversity. This contributes to improving the model's generalization performance and robustness. The SA attention strategy can re-weight features in response to enhance attention to important regions, adaptively fusing contextual information of different scales.

Another interesting finding is that when the training set and the validation set share the same settings, the model's performance is improved, as shown in Table IV, method A involves separate training and validation sets, while method B uses the same set for training and validation. From a subjective perspective, the model in this situation should not be robust. In reality, data transformation techniques alleviate the impact of this situation, and the performance improvement is attributed, to some extent, to the slightly increased training data.

TABLE III.     ELANET USES DROPOUT STRATEGIES FOR ABLATION STUDIES

| Dp | At | P | R | mAP@0.5 | mAP@0.5-0.95 |
|---|---|---|---|---|---|
| -- | -- | 0.855 | 0.810 | 0.870 | 0.463 |
| √ | -- | 0.848 | 0.790 | 0.855 | 0.457 |
| -- | √ | 0.850 | 0.803 | 0.868 | 0.472 |
| √ | √ | 0.882 | 0.816 | 0.887 | 0.501 |

TABLE IV.     COMPARISON OF ELANET PERFORMANCE WITH DIFFERENT DATASET SETTINGS

| Method | P | R | mAP@0.5 | mAP@0.5-0.95 |
|---|---|---|---|---|
| A | 0.872 | 0.779 | 0.864 | 0.468 |
| B | 0.882 | 0.816 | 0.887 | 0.501 |

## IV.    DISCUSSION AND FUTURE WORK

In this study, the challenges of cross-domain object detection were explored, focusing on addressing distribution differences between different datasets. The proposed ELANet model, based on feature extraction and fusion, was introduced. Experimental results indicate that ELANet performs exceptionally well on the RSOD dataset, while its performance on the GWHD dataset is relatively lower. This difference may be attributed to the clearer features in the airplane images of the RSOD dataset, making it easier for the model to learn effective feature representations. In contrast, the complex background and the similarity in color between target objects and the background in the GWHD dataset increase the difficulty of model recognition.

Furthermore, this study compared ELANet's performance on the RSOD and GWHD datasets with other methods,

revealing superior performance on both datasets. This suggests that ELANet can better handle cross-domain challenges and improve the accuracy of object detection. Nevertheless, there are still some issues to be addressed. For example, ELANet's performance is influenced by data preprocessing and training details. To further enhance the model's performance, deeper research into data preprocessing and training techniques is needed. Distribution differences have been a long-standing issue in machine learning [31]-[35], and it is hoped that this work will further stimulate researchers' interest in addressing this problem.

In future research, potential improvement avenues can be explored in the following directions:

Adversarial Robustness Learning: Conduct in-depth research to assess ELANet's robustness against adversarial attacks. Strengthening the model's security is crucial for real-world deployment and addressing potential security challenges.

Testing in Complex Environments: Test ELANet's robustness in more complex environmental conditions, especially in scenarios with natural factors such as varying lighting and rainy weather. Consider employing more powerful data augmentation and processing techniques to enhance the model's adaptability to these challenges.

Driving Agricultural Technology Innovation: Expand ELANet's application to more agricultural domains, such as agricultural robots, precision agriculture, and plant disease diagnosis. Through widespread application in agricultural technology, ELANet aims to provide more intelligent and efficient solutions for agricultural production.

## V.    CONCLUSION

In this study, the focus was on exploring a visual model for domain adaptation without the need for explicit adaptation, particularly addressing the challenge of domain adaptation in supervised graph learning tasks. In this work, the ELANet model was proposed, innovatively introducing the ELA module and integrating it into the feature decoder, successfully achieving high-quality domain adaptation in the remote sensing image domain. The model demonstrated outstanding performance on real datasets. The research indicates that regularization techniques are crucial for mitigating the domain variance between training and testing data, especially when incorporating a module that reweights features in response. Importantly, the use of the ELANet model not only improved accuracy but also enhanced efficiency. In experiments validating the generality and domain adaptation of ELANet, cross-domain settings were employed, demonstrating the model's generality and domain adaptability. Overall, the introduction of ELANet is expected to advance research in domain adaptation, providing an innovative approach to addressing domain adaptation challenges in supervised graph learning tasks.

REFERENCES

[1]  K. Li, G. Wan, G. Cheng, L. Meng, and J. Han, "Object detection in optical remote sensing images: A survey and a new benchmark," ISPRS Journal of Photogrammetry and Remote Sensing, vol. 159, January 2020, pp. 296-307.

[2] S. J. Pan and Q. Yang, "A Survey on Transfer Learning," IEEE Transactions on Knowledge and Data Engineering, vol. 22, no. 10, pp. 1345-1359, Oct. 2010.

[3] P. Dollar, C. Wojek, B. Schiele, and P. Perona, "Pedestrian Detection: An Evaluation of the State of the Art," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 34, no. 4, pp. 743-761, April 2012.

[4] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in 2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, 2012, pp. 2066-2073.

[5] Y. Ganin, E. Ustinova, H. Ajakan, et al., "Domain-adversarial training of neural networks," Journal of Machine Learning Research, vol. 17, no. 59, 2016, pp. 1-35.

[6] C. Doersch and A. Zisserman, "Multi-task Self-Supervised Visual Learning," in 2017 IEEE International Conference on Computer Vision (ICCV), Venice, Italy, 2017, pp. 2070-2079.

[7] J. Casebeer, N. J. Bryan, and P. Smaragdis, "Meta-AF: Meta-Learning for Adaptive Filters," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 31, 2023, pp. 355–370.

[8] W. Wang, J. Dai, Z. Chen, et al., "Internimage: Exploring large-scale vision foundation models with deformable convolutions," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023, pp. 14408-14419.

[9] P. Goel and A. Ganatra, "Unsupervised Domain Adaptation for Image Classification and Object Detection Using Guided Transfer Learning Approach and JS Divergence," Sensors, vol. 23, no. 9, 2023, pp. 4436.

[10] A. Elshamli, G. W. Taylor, A. Berg, et al., "Domain adaptation using representation learning for the classification of remote sensing images," IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, vol. 10, no. 9, 2017, pp. 4198-4209.

[11] J. Zhang, S. Xu, J. Sun, et al., "Unsupervised Adversarial Domain Adaptation for Agricultural Land Extraction of Remote Sensing Images," Remote Sensing, vol. 14, no. 24, 2022, pp. 6298.

[12] M. Valerio Giuffrida, A. Dobrescu, P. Doerner, et al., "Leaf counting without annotations using adversarial unsupervised domain adaptation," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, 2019, pp. 0-0.

[13] S. Marino, P. Beauseroy, A. Smolarz, "Unsupervised adversarial deep domain adaptation method for potato defects classification," Computers and Electronics in Agriculture, vol. 174, 2020, 105501.

[14] A. Fuentes, S. Yoon, T. Kim, et al., "Open set self and across domain adaptation for tomato disease recognition with deep learning techniques," Frontiers in Plant Science, vol. 12, 2021, 758027.

[15] X. Wu, X. Fan, P. Luo, et al., "From Laboratory to Field: Unsupervised Domain Adaptation for Plant Disease Recognition in the Wild," Plant Phenomics, vol. 5, 2023, 0038.

[16] F. Magistri, J. Weyler, D. Gogoll, et al., "From one field to another— Unsupervised domain adaptation for semantic segmentation in agricultural robotics," Computers and Electronics in Agriculture, vol. 212, 2023, 108114.

[17] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition," in 2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops, Columbus, OH, USA, 2014, pp. 512-519.

[18] H. Lu, C. Shen, Z. Cao, Y. Xiao, and A. van den Hengel, "An Embarrassingly Simple Approach to Visual Domain Adaptation," IEEE Transactions on Image Processing, vol. 27, no. 7, 2018, pp. 3403–3417.

[19] E. David, M. Serouart, D. Smith, et al., "Global Wheat Head Detection 2021: An Improved Dataset for Benchmarking Wheat Head Detection Methods," Plant Phenomics, Sep 22, 2021;2021:9846158.

[20] Y. Long, Y. Gong, Z. Xiao, and Q. Liu, "Accurate Object Localization in Remote Sensing Images Based on Convolutional Neural Networks," in IEEE Transactions on Geoscience and Remote Sensing, vol. 55, no. 5, May 2017, pp. 2486-2498.

[21] H. Zhu, X. Chen, W. Dai, K. Fu, Q. Ye, and J. Jiao, "Orientation Robust Object Detection in Aerial Images Using Deep Convolutional Neural Network," in 2015 IEEE International Conference on Image Processing (ICIP), Quebec City, QC, Canada, 2015, pp. 3735-3739.

[22] C. Wang, H. Liao, and I. Yeh, "Designing Network Design Strategies Through Gradient Path Analysis," arXiv preprint arXiv:2211.04800, 2022.

[23] G. Jocher, A. Stoken, A. Chaurasia, J. Borovec, Y. Kwon, K. Michael, et al., "Yolov5 Repository," Available at https://github.com/ultralytics/yolov5, Accessed on 10/28/2023.

[24] W. Lv, S. Xu, Y. Zhao, et al., "Detrs beat yolos on real-time object detection," arXiv preprint arXiv:2304.08069, 2023.

[25] G. Yang, J. Lei, Z. Zhu, et al., "AFPN: Asymptotic Feature Pyramid Network for Object Detection," arXiv preprint arXiv:2306.15988, 2023.

[26] Q. L. Zhang and Y. B. Yang, "SA-Net: Shuffle Attention for Deep Convolutional Neural Networks," in ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2021, pp. 2235-2239.

[27] A. Paszke, S. Gross, F. Massa, et al., "PyTorch: An Imperative Style, High-Performance Deep Learning Library," Advances in Neural Information Processing Systems, vol. 32, 2019.

[28] D.P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," arXiv preprint arXiv:1412.6980, 2014.

[29] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 39, no. 6, pp. 1137-1149, June 1, 2017.

[30] C. Y. Wang, A. Bochkovskiy, and H. Y. M. Liao, "YOLOv7: Trainable Bag-of-Freebies Sets New State-of-the-Art for Real-Time Object Detectors," in 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Vancouver, BC, Canada, 2023, pp. 7464-7475.

[31] Q. Ming, L. Miao, Z. Zhou, and Y. Dong, "CFC-Net: A Critical Feature Capturing Network for Arbitrary-Oriented Object Detection in Remote-Sensing Images," in IEEE Transactions on Geoscience and Remote Sensing, vol. 60, pp. 1-14, 2022, Art no. 5605814.

[32] X. Hu and C. Zhu, "Shared-Weight-Based Multi-Dimensional Feature Alignment Network for Oriented Object Detection in Remote Sensing Imagery," Sensors, vol. 23, no. 1, 2022, Article 207.

[33] Q. Ming, L. Miao, Z. Zhou, J. Song, and X. Yang, "Sparse Label Assignment for Oriented Object Detection in Aerial Images," Remote Sensing, vol. 13, no. 14, 2021, Article 2664.

[34] S. Falahat and A. Karami, "Maize Tassel Detection and Counting Using a YOLOv5-Based Model," Multimedia Tools and Applications, vol. 82, pp. 19521–19538, 2023.

[35] Y. Yang and S. Soatto, "FDA: Fourier Domain Adaptation for Semantic Segmentation," in 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Seattle, WA, USA, 2020, pp. 4084-4094.

# Improving Brain Tumor MRI Image Classification Prediction based on Fine-tuned MobileNet

Quy Thanh Lu, Triet Minh Nguyen, Huan Le Lam

Information Technology Department, FPT University, Can Tho, Viet Nam

*Abstract*—Brain tumors are a prevalent issue in contemporary society as they impact human health. The location of the tumor in the brain determines the variety of symptoms that may manifest. Some frequent symptoms are cephalalgia, convulsions, visual impairments, nausea, emesis, asthenia, paresthesia, dysphasia, personality alterations, and amnesia. The prognosis for brain cancer differs considerably depending on the cancer type. Nevertheless, brain tumors are amenable to treatment with surgical intervention, chemotherapy, and radiotherapy if the diagnosis is timely. Furthermore, artificial intelligence and machine learning can assist in the detection of brain tumors as they have significant implications for the analysis of Magnetic Resonance Imaging (MRI). To accomplish this objective, automated measurement instruments were proposed based on the processing of MRI. In this study, we employed the latest developments in deep transfer learning and fine-tuning to identify tumors without many complex steps. We gathered data from authentic MRI of 3264 subjects (i.e., 926 glioma tumors, 937 meningioma tumors, 901 pituitary tumors, and 500 normal). With the MobileNet model from the Keras library, we attained the highest validation accuracy, test accuracy, and F1 score in four-class classifications was 97.24%, 97,86%, and 97.85%, respectively. Concerning two-class classification, high accuracy values were obtained for most of the models (i.e., ~100%). These outcomes and other performance indicators demonstrate a strong capability to diagnose brain tumors from conventional MRI. The current research developed a supportive machine learning that can aid doctors in making the accurate diagnosis with less time and mistakes.

*Keywords—Brain tumor; fine-tuning; transfer learning; Magnetic Resonance Imaging (MRI); MobileNet*

## I. INTRODUCTION

The health of people in the modern world is adversely affected by numerous diseases, among which brain tumors are prevalent in various age groups, from adolescents to seniors. Brain tumors are not a helpless disease and there are several treatment options available for patients with this disease. Some of the most common and efficacious methods of treating brain tumors are surgery, chemotherapy, radiation therapy, or a combination of these techniques. Surgery entails excising the tumor or part of it through an operation, while chemotherapy and radiation therapy employ drugs and high-energy rays to eradicate or reduce the tumor cells. The selection of treatment for each patient hinges on multiple factors, such as the type of the tumor, the grade of the tumor, the location of the tumor, the size of the tumor, as well as the age and general health of the patient. By considering these factors, physicians can devise the optimal treatment plan for each patient and enhance their prospects of recovery.

Nevertheless, it is still dangerous if we avoid it. For example, A brain tumor is a grave condition irrespective of its benignity or malignancy. Tumors within the cranium expand and exert pressure on different regions of the brain, impairing their function. A glioma is a tumor that emerges when glial cells proliferate abnormally. Typically, these cells support neurons and facilitate the operation of your central nervous system. Gliomas commonly grow in the brain, but can also emerge in the spinal cord. Gliomas are malignant (cancerous), but some can be very slow-growing. In addition, Meningioma is the most frequent type of primary brain tumor, it is a tumor that develops from the meninges, the protective layers around the brain and spinal cord. Meningiomas occur more often in women and are usually detected at older ages. Moreover, Pituitary tumors are abnormal enlargements that originate in the pituitary gland. Some of these tumors cause the pituitary gland to produce an excess of certain hormones that regulate vital body functions. tumors stimulate the adrenal glands to produce too much cortisol. This causes a condition called Cushing disease. Others can cause the pituitary gland to produce too little of those hormones.

Numerous data gathered in recent years globally indicate the detrimental effects of brain tumors. The CBTRUS Statistical Report states that from 2016 to 2020, 86,030 deaths occurred due to malignant brains. This corresponds to an average annual mortality rate of 4.42 and an average of 17,206 deaths per year [1]. In 2030, 145,650 new cases of brain tumors are projected to be diagnosed in China with 68,730 men and 76,920 women [2]. The 5-year relative survival rate after diagnosis of a malignant brain tumor was 35.7% and for a non-malignant brain tumor was 91.9% [1]. The most frequent malignant brain tumor was glioblastoma with 14.2% of all tumors, and the most common non-malignant tumor was meningioma with 40.8% of all tumors [1]. Glioblastoma was more prevalent in males, and meningioma was more prevalent in females. In children and adolescents aged 0-19 years, the incidence rate of primary tumors was 6.14 out of 100,000 people [3]. An estimated 3,920 new cases of primary childhood brain tumors are expected to be diagnosed in 2023 [4]

Recently, Magnetic Resonance Imaging (MRI) has been considered one of the most efficient methods for identifying the irregular parts of the central nervous system and human brain. Moreover, the early diagnosis of brain tumors is one of the crucial tasks and offers many advantages to patients. Early identification of tumors helps doctors devise suitable treatment plans and helps lower mortality in patients with brain tumors as quickly as possible. There are many methods that clinicians can use to detect brain tumors as Computerized Tomography

(CT) or MRI. However, it would be time-intensive to diagnose using the MRI method. Furthermore, doctors have to prepare and conduct many procedures to finish a standardized process for patients. Therefore, applying the newest advanced techniques in the diagnostic process can spare valuable time for doctors. Moreover, it can provide doctors with a recommendation to improve the diagnostic process and increase the outcomes.

Artificial intelligence is one of the most prominent advanced techniques invented in recent years. Therefore, we decided to use it to assist doctors in diagnosing MRI images. In the field of artificial intelligence, machine learning is related to the development and research of statistical algorithms that can effectively generalize data and syntax. Based on the trained and prepared results, the computer will perform tasks without explicit instructions [5]. Machine learning algorithms have many applications. For example, financial prediction, transportation, education, data structure in health care systems, drug reaction prediction, diabetes research, cyber security, banking and finance, and social media [6].

Our study used transfer learning and fine-tuning, a part of deep learning which is the subset of the machine learning method. That aim enables us to reuse pre-trained models for new tasks and datasets. Moreover, Transfer learning concentrates on transferring general knowledge from one domain to freeze certain layers to preserve general [7]. Fine-tuning adapts the model to a particular allows the pre-trained layers to be updated [8]. We suggest a new method to use the MobileNet model in the Keras library in a Convolutional Neural Network (CNN) which is appropriate for image recognition and processing tasks and it has achieved state-of-the-art outcomes on a wide range of image recognition tasks, such as object classification, object detection, and image segmentation [9]. It is trained using a huge dataset of annotated images. Once trained, a CNN can be utilized to classify new images or extract features for use in other applications such as object detection or image segmentation.

The contributions of this paper are as follows:

- We propose a complete reliable artificial intelligence model that is used for brain tumor classification including glioma, meningioma, and pituitary tumor. Hence, it can create an easy and speedy way for doctors to detect and classification on their medical purpose.

- Our method achieves a high performance (i.e., <97%) of the deep learning models for four-class (glioma, meningioma, pituitary, and normal) and pair-wise classification problems also achieve a high success (i.e., ~100%). As a result. The effectiveness of each model in terms of training and testing times was also evaluated.

- Our collected MRIs of subjects afflicted with brain tumors, as well as healthy ones, as verified by the specialists in the hospital. This dataset is confirmed for the development of automated machine learning and AI algorithms for the detection of brain tumors and can be applied to educating medical students.

- We point out that GradCam can be used effectively for visual explanations. So, highlighting the important regions (i.e., glioma tumor, meningioma tumor, and pituitary tumor) in the image can provide a more intuitive view for physicians

Our research paper comprises four main sections. In the subsequent Section II, we indicate some of the related research that we employed for references. Following the related research section is the methodology Section III, this section elucidates in detail all of the methods utilized in the article. Subsequently, Section IV will refer to the experiments, and how we conduct and evaluate the accuracy of the deep learning model. Lastly, in the final Section V, we summarize our article and examine the essential domains associated with the study.

## II. RELATED WORK

Besides the change in environment and population, a lot of diseases have negative effects on human life and one of these is brain tumors. The survey pointed out that 1,323,121 people living with brain and other CNS tumors (malignant and non-malignant) on December 31, 2019 [1]. At that time, a bunch of research on medical and artificial intelligence helped humans in the healing of those sicknesses. Wadhah Ayadi detected brain tumors by suggesting a new CNN that contains various layers such as convolution, Rectified Linear Unit (Relu), and pooling to achieve a best Accuracy is 94.74 [10]. Following Ahmad Saleh, His scholarly investigation endeavors to enhance the proficiency of MRI apparatus in the categorization of cerebral neoplasms and discerning their respective classifications, using AI Algorithms, CNN, and Deep Learning. The MobileNet model is used in this study. The evaluation of the image dataset is conducted utilizing the F1-score metric, yielding a commendable accuracy rate of 97.25% [11].

Machine Learning helped experts and doctors research medical image analysis. This led to, a change in normal examination and treatment to aid medical procedures to avoid a waste of time and money. Hence, typical studies appear more and more such as Wadhah Ayadi. The presented methodology used normalization, dense speeded-up robust features, and histogram of gradient techniques to enhance the quality of MRI and produce a discriminative feature set. The accuracy attained through the implementation of this method is measured at 90.27% [12]. In addition, Muhammad Imran Sharif suggests Densenet201 Pre-Trained Deep Learning Model. The attributes of the trained model are derived from the average pooling layer, elucidating the profound information on each specific tumor type. In addition, He includes two new models Entropy–Kurtosis-based High Feature Values (EKbHFV) and modified genetic algorithm (MGA). Finally, the research paper has achieved an accuracy higher than 95% [13].

Convolutional Neural Network (CNN) is one of the parts of deep learning models commonly used in Computer Vision that help computers understand and interpret images or visual data. It has main contributions to creating intelligent systems with great accuracy. So, it has produced several academic works as S Kumar's research paper. In short, the technical use of the Deep Convolution Neural Network (Deep CNN) for performing the brain tumor classification with Dolphin-SCA as the training algorithm. The database is MRI images given by

the BRATS database and SimBRATS, and the suggested model has shown a maximum accuracy of 96.39% [14]. Moreover, Díaz-Pernas applied the method on a publicly available MRI image dataset of 3064 images from 233 patients compared with previously classical segmentation and classification published methods. In this comparative analysis, the proposed method demonstrated exceptional results, achieving an impressively high tumor classification accuracy of 0.973 [15].

Processing medical images by automatic segmentation and classification becoming extremely important around the world [16], [17], especially in the medical field such as diagnostics, growth prediction, and treatment of brain tumors. As a result, a patient can save their life because an early detection of brain tumors that helps to increase their survival rate. Applying machine learning, the brain classification paper from Huong Hoang Luong points out that the K-mean clustering algorithm stratifies the samples into three distinct view angles of MRI, namely, transverse, coronal, and sagittal planes. This strategic classification process was coupled with the integration of a modified Residual Network (ResNet) architecture. Finally, He reached a performance of 96% in the brain tumor classification accuracy [18]. Furthermore, with another k-means algorithm S. Rinesh's innovative approach outperformed conventional methods such as hybrid k-means clustering and parallel k-means clustering, exhibiting superior results with a higher peak signal-to-noise ratio and a reduced mean absolute error value. The proposed model attained an accuracy of 96.47% [19].

Additionally, we have introduced an automated classification system designed for the intricate challenge of categorizing multiclass brain tumor MRI. This task, inherently more complex and demanding, transcends the relative simplicity of binary classification. The dataset is almost equal to Khan Swatithe's paper uses a pre-trained deep CNN model and introduces a block-wise fine-tuning strategy rooted in transfer learning principles. Notably, this methodology is characterized by its generality, eschewing the need for handcrafted features, and demanding minimal preprocessing. Impressively, the proposed approach attains an average accuracy rate of 94.82% within the context of a five-fold cross-validation framework [20]. Besides, with a straightforward architectural design and the absence of any antecedent region-based segmentation, Nyoman Abiwinanda achieved commendable results, attaining a training accuracy of 98.51% and a peak validation accuracy of 84.19%. Notably, these outcomes stand in favorable comparison to the performance of more intricate region-based segmentation algorithms [21].

In summary, the related work saw notable weaknesses, particularly characterized by low accuracy. Many current models struggle to consistently achieve high precision across diverse datasets, leading to concerns about their robustness and generalizability. Additionally, the lack of a visual explanation such as GradCam for evaluating and comparing models hinders progress. Addressing these limitations is imperative to propel deep learning toward more reliable and universally applicable solutions, ensuring advancements that transcend the current constraints of accuracy and evaluation methodologies.

## III. METHODOLOGY

### A. The Research Implementation Procedure

In this research paper, we propose a method including 11 steps from input to output shown in Fig. 1. The roles of the steps are shown as follows:



Fig. 1. The implementing procedure flowchart.

*1) Collecting dataset:* The dataset is collected by Swati Kanchan in B. Tech Dept of CSE from NIT Durgapur. The dataset comprises MRI images, including three types of brain tumors—meningioma, glioma, and pituitary—as well as normal images. This comprehensive collection serves as a valuable resource for medical research and diagnostic advancements.

*2) Pre-processing image:* This important step requires both resizing and normalization to establish standardized input conditions for machine-learning models. This leads to an increase in the outcome of the results.

*3) Dividing* the dataset into three categories train validation and test: In total, 3264 images constitute the comprehensive MRI dataset, with random selection for training, validation, and testing phases. The datasets are randomly chosen using an 8-1-1 scale, allocating 8 parts for training, 1 for validation, and 1 for testing, ensuring a balanced distribution for robust model development and evaluation.

*4) Dividing the training set into folders:* Types of brain tumors (i.e., meningioma, glioma, and pituitary and normal) are divided into many different folders. The biggest include 4 classes (i.e., meningioma, glioma, and pituitary and normal). We want to show the training data more exactly way and use it to compare the biggest folder. Dense, the other folder contains two class classifications includes: meningioma-normal, glioma-normal, and pituitary-normal

*5) Building the model:* To conduct experiments, we reconstructed the model based on the prototype of the CNN

architecture by inheriting its core processing layers and modifying some layers to achieve improved results. As a result, the model creates an excellent result for our training test with Keras's models library.

*6) Applying transfer learning:* In transfer learning, the model gets trained on a big set of data for a specific job. This dataset could be general or have lots of labeled information. The things the model learns are saved in its weights, especially in the lower layers. These layers gather important features from the input data, making it easier for the model to understand and work with the information.

*7) Validating and collecting the accuracy score:* After completing the training of the model, we assessed its precision by summarizing the training accuracy derived from the model's predictions. Subsequently, we evaluated the accuracy of the test set using the initially separated testing set."

*8) Applying fine tuning:* In this procedure, the process of adjusting the hyperparameters of a model to improve its performance on a dataset and using the weights of an already trained model as the starting values for training a new model. Hyperparameters are the parameters that control the learning process and we finish after the model has been trained on a training set, and done when using a validation set.

*9) Validating, collecting and drawing results with GradCam:* We summary all the metrics like validation accuracy, test accuracy, and F1 score, then. Images used GradCam to represent detected brain tumors in color to more accurately represent the results.

*10)Reconstructing and comparing the cycles with other models:* When we have results in one model, we rework another model in Keras including MobileNet, Inception V3, VGG16, ResNet50, and EfficientNet B3 to compare the final result

*11)Showing the result:* After conducting a comparison, results will be presented through tables and graphs to facilitate meaningful comparisons.

### B. Pre-processing Image

In the pre-processing pipeline for images, a pivotal step involves both resizing and normalization to establish standardized input conditions for machine learning models. The resizing operation transforms the input image I to a uniform dimension of (new width, new height). As a result, we decided to choose 224 pixels for weight and 224 pixels for height (1) using a chosen interpolation method, as captured by the formula:

$$I_{resize}(new_{width}, new_{height}) = Resize(I, 244,244) \quad (1)$$

This operation (1) is essential for establishing a consistent input size, a prerequisite in various deep-learning applications. Following resizing in equation (2), the subsequent normalization process adjusts pixel values to a range between 0 and 1, facilitating model training and convergence. The normalization is mathematically expressed as:

$$O(x,y) = \frac{I_{resize}(x,y)}{255} \quad (2)$$

In this Eq. 2, $O(x,y)$ signifies the output pixel value at position $(x,y)$ in the preprocessed image. The division by 255 ensures that the pixel values are scaled to fit within the [0, 1] range, aligning with common conventions in image processing.

This dual procedure of resizing to (224, 224) (1) and subsequent normalization not only ensures a standardized size for all images but also provides a consistent pixel value scale (2), thereby enhancing the efficiency and robustness of downstream machine learning tasks.

### C. Transfer Learning and Fine Tuning of CNN (MobileNet)

We noticed that transfer learning and fine-tuning of CNN have become pivotal techniques in the realm of computer vision, enabling the effective utilization of pre-trained models on new tasks [22][23][24]. One such exemplary model is MobileNet, a lightweight and efficient architecture design. MobileNet, characterized by depth wise separable convolutions, significantly reduces the computational burden while preserving the model's capacity to capture intricate features in images.

The MobileNet model's architecture is inherently rooted in the principles of convolutional neural networks, with its convolutional layers serving as feature extractors. This architecture excels in tasks demanding real-time performance and resource efficiency. When considering transfer learning, MobileNet offers a valuable starting point. Pre-trained on large-scale image datasets, it possesses a robust understanding of general image features, making it an ideal candidate for various computer vision applications.

The integration of MobileNet into the broader CNN architecture is seamless and complementary such as Fig. 2. The initial layers of a CNN, responsible for low-level feature extraction, can be substituted with the MobileNet backbone. This strategic replacement allows the model to retain its ability to recognize high-level features while benefiting from MobileNet's efficiency in processing low-level features. The resulting hybrid architecture strikes a balance between computational efficiency and task-specific adaptability.

In conclusion, the fusion of transfer learning and fine-tuning within the realm of CNNs, particularly with the utilization of the MobileNet model, represents a powerful approach to solving diverse computer vision challenges. This leads to the final results about validation accuracy, test accuracy, and F1 score increase and reaching a desirable point.

### D. Visual Explanation by GradCam

In our paper, we decided to choose a visual explanation by Gradient-weighted Class Activation Mapping (i.e., GradCam), because it stands as a pivotal technique in unraveling the decision-making processes of CNN. It operates by shedding light on crucial regions within an image that heavily influence the model's final prediction, thus enhancing both interpretability and confidence in the outputs.

Fig. 2. Procedure of transfer learning and fine-tuning in ours model with custom layers.

In more detail, GradCam relies on the gradient information flowing into the final convolutional layer of a CNN. Let $F^k$ (3) represent the k-th feature map from the last convolutional layer, and $w_c^k$ (3) show the weight of the k-th feature map for the target class 'c.' The GradCam heatmap $L_c^{GradCam}$ (3) is computed as the global average pooling of the positive gradients:

$$L_c^{GradCam} = ReLU(\sum_k w_c^k \cdot F^k) \qquad (3)$$

The ReLU (3) function ensures that only positive contributions are considered, highlighting the regions with a positive influence on the decision for the target class 'c' Fig. 3.

To explain this formula in detail, let's express the weight $w_c^k$ (4) for a specific feature map 'k' as the global average pooling of the gradients:

$$w_c^k = \frac{1}{Z} \sum_{i,j} \frac{\partial y_c}{\partial A_{i,j}^k} \qquad (4)$$



Fig. 3. GradCam functionality.

Because $y_c$ (4) denotes the logit for the target class 'c,' and $A_j^k$ (4) presents the activation of the k-th feature map at position (i, j). The normalization term 'Z' ensures that the weights sum to 1, providing a meaningful contribution measure.

The computation of the gradient $\frac{\partial y_c}{\partial A_j^k}$ (5) involves back-propagating the derivative of the logit with respect to the activation of the k-th feature map at position (i, j). By means of mathematics, this can be expressed as:

$$\frac{\partial y_c}{\partial A_j^k} = \frac{\partial y_c}{\partial F_{i,j}^k} \cdot \frac{\partial F_{i,j}^k}{\partial A_{i,j}^k} \qquad (5)$$

This chain of derivatives involves the gradient of the logit with respect to the activation of the k-th feature map at position (i, j), $\frac{\partial y_c}{\partial F_{i,j}^k}$ (5), and the gradient of the activation at (i, j) with respect to the input activation, $\frac{\partial F_{i,j}^k}{\partial A_{i,j}^k}$ (5).



Fig. 4. The final result after applying a visual explanation by GradCam.

By using this GrandCam, our results have the contribution of each feature map to the final decision in Fig. 4, and facilitate future use by professionals and doctors.

## IV. EXPERIMENTS

### A. Dataset and Peformance Metrics

The research used a single dataset for both the training, validation, and testing phases in this analysis. The data was taken and augmented by Swati Kanchan in B. Tech Dept of CSE from NIT Durgapur, 3264 images constitute the comprehensive MRI dataset in total including 926 glioma tumors, 937 meningioma tumors, 901 pituitary tumors, and 500 no tumor.

Moreover, the performance of the models was assessed using five metrics: validation accuracy, test accuracy, precision, recall, and F1 score play pivotal roles in assessing the performance and generalization capabilities of a trained model.

Validation accuracy (Val acc) in Eq. 6 represents the model's precision on a separate dataset during training, measuring its ability to learn without overfitting the training data. By means of mathematics, validation accuracy is computed as:

$$Val\ acc = \frac{Total\ Number\ of\ instance\ in\ validation\ set}{Number\ of\ correctly\ predictec\ instrances} \quad (6)$$

Test accuracy (Test acc) in Eq. 7 reflects the model's proficiency in making accurate predictions on previously unseen data, providing insights into its real-world applicability. This metric is calculated by:

$$Test\ acc = \frac{Total\ number\ of\ instances\ in\ test\ set}{Number\ of\ correctly\ predicted\ instances} \quad (7)$$

Recall in Eq. 8, a metric crucial in scenarios where identifying true positives is paramount, is defined as:

$$Recall = \frac{True\ Positives}{True\ Positives + False\ Negatives} \quad (8)$$

Precision in Eq. 9 is a fundamental metric in the evaluation of classification models. Mathematically, precision is defined as the ratio of true positives (instances correctly predicted as positive) to the sum of true positives and false positives (instances incorrectly predicted as positive). The precision formula is given by:

$$Precision = \frac{True\ Positives}{True\ Positives + False\ Positives} \quad (9)$$

The F1 score in Eq. 10, a harmonic mean of precision and recall, is expressed as:

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (10)$$

### B. Scenario 1: The Results of Classifying MRI Images Into Two Classes: Normal or Glioma Tumor

The purpose of the experiments was to evaluate the effectiveness of the pre-trained models, after customization and training, in identifying the correct disease diagnosis of the MRI image. Moreover, these data can help us make easier and more intuitive comparisons across brain tumor categories (four classes or two classes).

Tables I and II show the performance evaluation metrics for classifying MRI images into normal or scoliosis. In transfer learning, the Resnet-50 achieved the highest accuracy (i.e., 100%). After fine-tuning, three models including MobileNet, InceptionV3, and EfficientNetB3 produced results that exceeded expectations (i.e., 100% for three models). In contrast, Resnet-50 showed worse results than the previous experiment.

TABLE I. THE ACCURACY OF CLASSIFYING MRI IMAGES INTO TWO CLASSES: GLIOMA TUMOR AND NORMAL IN TRANSFER LEARNING, FOR EACH DEEP LEARNING MODEL

| Model | Transfer learning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 99.30% | 98.60% | 98.66% | 98.60% | 98.61% |
| ResNet50 | 100.00% | 98.60% | 98.60% | 98.60% | 98.60% |
| VGG16 | 99.30% | 97.90% | 98.02% | 97.90% | 97.92% |
| **Ours** | **99.30%** | **97.90%** | **97.92%** | **97.90%** | **97.91%** |
| InceptionV3 | 99.30% | 93.01% | 93.11% | 93.01% | 93.04% |

TABLE II. THE ACCURACY OF CLASSIFYING MRI IMAGES INTO TWO CLASSES: GLIOMA TUMOR AND NORMAL IN FINE TUNING, FOR EACH DEEP LEARNING MODEL

| Model | Fine tuning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 100.00% | 99.30% | 99.31% | 99.30% | 99.30% |
| ResNet50 | 99.30% | 96.50% | 96.64% | 96.50% | 96.53% |
| VGG16 | 99.30% | 87.41% | 88.33% | 87.41% | 86.83% |
| **Ours** | **100.00%** | **97.90%** | **97.92%** | **97.90%** | **97.91%** |
| InceptionV3 | 100.00% | 98.60% | 98.60% | 98.60% | 98.60% |

Fig. 5 and Fig. 6 show a sample training and validation progress curve showing the loss and accuracy values of ours in fine-tuning. The figure displays a stable learning behavior and appropriate training and validation sets.



Fig. 5. Training accuracy and validation accuracy in fine-tuning of ours model (Glioma tumors and normal).

Fig. 6.    Training loss in and validation loss in fine-tuning of ours model (Glioma tumors and normal).

Fig. 7 shows the confusion matrix for a sample run of ours model in two classes. In that run, the number of testing images is 143.



Fig. 7.    Confusion matrix in fine tuning for ours model (Glioma tumors and normal).

## C. Scenario 2: The Results of Classifying MRI Images Into Two Classes: Normal or Meningioma Tumor

Tables III and IV showed similar results in scenario 2 with the Resnet-50 achieving the biggest accuracy (i.e., 100%) in transfer learning. After fine-tuning, ours achieved extremely impressive results with low test loss and all other aspects reaching 100%.

TABLE III.    THE ACCURACY OF CLASSIFYING MRI IMAGES INTO TWO CLASSES:  MENINGIOMA TUMOR AND NORMAL IN TRANSFER LEARNING, FOR EACH DEEP LEARNING MODEL

| Model | Transfer learning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 99.31% | 98.61% | 98.64% | 98.61% | 98.60% |
| ResNet50 | 100.00% | 96.53% | 96.84% | 96.53% | 96.56% |
| VGG16 | 97.92% | 97.92% | 97.92% | 97.92% | 97.91% |
| **Ours** | **97.22%** | **95.14%** | **95.17%** | **95.14%** | **95.15%** |
| InceptionV3 | 95.83% | 97.22% | 97.43% | 97.22% | 97.25% |

TABLE IV.    THE ACCURACY OF CLASSIFYING MRI IMAGES INTO TWO CLASSES: MENINGIOMA TUMOR AND NORMAL IN FINE TUNING, FOR EACH DEEP LEARNING MODEL

| Model | Fine tuning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 100.00% | 99.31% | 99.32% | 99.31% | 99.31% |
| ResNet50 | 100.00% | 86.11% | 90.08% | 86.11% | 86.44% |
| VGG16 | 99.31% | 95.14% | 95.17% | 95.14% | 95.15% |
| **Ours** | **100.00%** | **100.00%** | **100.00%** | **100.00%** | **100.00%** |
| InceptionV3 | 100.00% | 97.22% | 97.22% | 97.22% | 97.22% |

In Fig. 8 and Fig. 9, we can see the detail of the training accuracy and training loss about two classes that is pituitary tumors and normal MRI images. This chart can show us the high result of a classification of brain tumors when using ours model.



Fig. 8.    Training accuracy and validation accuracy in fine-tuning of ours model (Meningioma tumors and normal).



Fig. 9.    Training loss in and validation loss in fine-tuning of ours model (Meningioma tumors and normal).

The confusion matrix in Fig. 10 shows that the result of model has a high performance when using it to classify brain tumor.

Fig. 10. Confusion matrix in fine tuning for ours model (Meningioma tumors and normal).

### D. Scenario 3: The Results of Classifying MRI Images Into Two Classes: Normal or Pituitary Tumor

In scenario 3, Tables V and VI show that ResNet-50 has top results in all aspects (i.e., 100%) while it decreased a litter bit after fine-tuning, other models have a bit of a transformation (i.e., ~99%) when finished two steps.

TABLE V. THE ACCURACY OF CLASSIFYING MRI IMAGES INTO TWO CLASSES: PITUITARY TUMOR AND NORMAL IN TRANSFER LEARNING, FOR EACH DEEP LEARNING MODEL

| Model | Transfer learning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 100.00% | 98.58% | 98.58% | 98.58% | 98.57% |
| ResNet50 | 100.00% | 100.00% | 100.00% | 100.00% | 100.00% |
| VGG16 | 94.29% | 98.58% | 98.58% | 98.58% | 98.58% |
| **Ours** | **100.00%** | **99.29%** | **99.30%** | **99.29%** | **99.29%** |
| InceptionV3 | 99.29% | 98.58% | 98.58% | 98.58% | 98.58% |

TABLE VI. THE ACCURACY OF CLASSIFYING MRI IMAGES INTO TWO CLASSES: PITUITARY TUMOR AND NORMAL IN FINE TUNING, FOR EACH DEEP LEARNING MODEL

| Model | Fine tuning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 100.00% | 99.29% | 99.30% | 99.29% | 99.29% |
| ResNet50 | 100.00% | 99.29% | 99.30% | 99.29% | 99.29% |
| VGG16 | 94.29% | 35.46% | 12.57% | 35.46% | 18.57% |
| **Ours** | **100.00%** | **99.29%** | **99.30%** | **99.29%** | **99.29%** |
| InceptionV3 | 100.00% | 99.29% | 99.30% | 99.29% | 99.29% |

In this experiment, Fig. 11 and Fig. 12 give an explanation of training accuracy and training loss in two classes of normal and pituitary tumors in the 'Ours' model.

Finally, the result confusion matrix is represented in Fig. 13 which shows that the performance of the 'Ours' model is very successful.



Fig. 11. Training accuracy and validation accuracy in fine-tuning of ours model (Pituitary tumors and normal).



Fig. 12. Training loss in and validation loss in fine-tuning of ours model (Pituitary tumors and normal).



Fig. 13. Confusion matrix in fine tuning for ours model (Pituitary tumors and normal).

### E. Scenario 4: The Results of Classifying MRI Images into Four Classes: Normal, Glioma Tumor, Meningioma Tumor, and Pituitary Tumor

Tables VII and VIII show the performance evaluation metrics for classifying MRI images into normal, glioma tumor, meningioma tumor, and pituitary tumor. The ResNet50 achieved the highest accuracy value in transfer learning over the three statistical measures with a validation accuracy of 94,17%, test accuracy of 92.05%, and F1 score of 92.06%. On

the other hand, ours model performed the best after fine-tuning with a validation accuracy of 97,24%, test accuracy of 97.86%, and F1 score of 97.86%. The other performance metrics display a consistent and homogenous ability to identify negative as well as positive cases with a similar performance pattern to the accuracy results (i.e., Ours model achieving the best results). The significance of the F1 score lies in its role as an evaluation metric specifically designed for classification problems. An F1 score serves as an indicator of the model's accuracy, emphasizing its ability to achieve both high precision and recall.

TABLE VII. THE ACCURACY OF CLASSIFYING MRI IMAGES INTO FOUR CLASSES: NORMAL, GLIOMA TUMOR, MENINGIOMA TUMOR, AND PITUITARY TUMOR IN TRANSFER LEARNING, FOR EACH DEEP LEARNING MODEL

| Model | Transfer learning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 93.55% | 91.13% | 91.50% | 91.13% | 91.16% |
| ResNet50 | 94.17% | 92.05% | 92.18% | 92.05% | 92.06% |
| VGG16 | 86.50% | 86.85% | 87.09% | 86.85% | 86.88% |
| **Ours** | **89.88%** | **85.93%** | **86.27%** | **85.93%** | **85.89%** |
| InceptionV3 | 83.13% | 81.65% | 81.79% | 81.65% | 81.69% |

TABLE VIII. THE ACCURACY OF CLASSIFYING MRI IMAGES INTO FOUR CLASSES: NORMAL, GLIOMA TUMOR, MENINGIOMA TUMOR, AND PITUITARY TUMOR IN FINE TUNING, FOR EACH DEEP LEARNING MODEL

| Model | Fine tuning | | | | |
|---|---|---|---|---|---|
| | Val acc | Test acc | Precision | Recall | F1 |
| EfficientNetB3 | 97.55% | 97.55% | 97.62% | 97.55% | 97.55% |
| ResNet50 | 97.24% | 95.11% | 95.20% | 95.11% | 95.09% |
| VGG16 | 65.03% | 15.29% | 2.34% | 15.29% | 4.06% |
| **Ours** | **97.24%** | **97.86%** | **97.91%** | **97.86%** | **97.86%** |
| InceptionV3 | 97.24% | 97.55% | 97.56% | 97.55% | 97.55% |

In Fig. 14 and Fig. 15 shows the training and validation progress curve for a sample run of the highest-performing model, which gives an indication of the fitting performance of the model and the need for more training. Training accuracy measures a model's performance on the training data, reflecting its ability to learn from the provided examples. Validation accuracy assesses the model's generalization to new, unseen data, helping identify potential overfitting or underfitting issues. Training loss quantifies the disparity between predicted and actual values in the training set, guiding the model to minimize errors during training. Validation loss mirrors this process on a separate dataset, serving as a key indicator of the model's generalization performance.

Fig. 16 shows the 'Ours' sample confusion matrix for four-class classification. This important step makes it possible for us to see a more intuitive comparison of the results achieved. Fig. 17 shows a sample output from the four-class classification process with the visual explanation by GradCam.



Fig. 14. Training accuracy and validation accuracy in fine-tuning of 'Ours' model (Four classes).



Fig. 15. Training loss in and validation loss in fine-tuning of 'Ours' model (Four classes).



Fig. 16. Confusion matrix in fine tuning for 'Ours' model (Four classes).

### F. Comparison with others State-of-the-art Methods

To examine the accuracy of the proposed model that our article has just given out in the previous section, we compare the accuracy score of the proposed model with other CNN architectures, which are VGG16, ResNet-50, ResNet-101 and DenseNet201.

Fig. 17. Output of four classes classification.

Accuracy, precision, recall, and the F1 score offer many aspects of performance. Accuracy gives overall correctness but may mislead in imbalanced datasets. Precision focuses on the accuracy of positive predictions. Recall assesses the model's ability to capture all relevant positives. The F1 score strikes a balance between precision and recall, making it invaluable for tasks with uneven class distributions. Comparing these metrics involves weighing trade-offs based on task-specific priorities. While accuracy provides a broad overview, precision and recall cater to nuanced aspects, and the F1 score harmonizes their interaction, ensuring a well-informed evaluation of classification models within the constraints of particular objectives. Finally, the result of getting the value of training and accuracy on the test set is illustrated as in Table IX.

TABLE IX. COMPARISON WITH OTHERS STATE-OF-THE-ART METHODS

| *Ref.* | *Proposed* | *Accuracy* |
|---|---|---|
| Wadhah Ayadi, et al. | CNN (ReLu) | 94.74% |
| Ahmad Saleh, et al. | MobileNet | 97.25% |
| Wadhah Ayadi, et al. | SVM | 90.27% |
| Muhammad Imran Sharif, et al. | Densenet201 | >95% |
| Sharan Kumar, et al. | Dolphin-SCA | 96.30% |
| Díaz-Pernas, et al. | Multiscale CNN | 97.30% |
| Huong Hoang Luong, et al. | ResNet-50 | 96% |
| Rinesh Sahadevan, et al. | MFNN | 96.46% |
| Khan Swati | VGG19 | 94.82% |
| Nyoman Abiwinanda | Multiscale CNN | 84.19% |
| **Proposed model** | | **97.86%** |

## V. CONCLUSION

In our project, we utilized transfer learning, a powerful technique in machine learning, to enhance our model's performance in identifying brain tumors from MRI images. Transfer learning involves leveraging knowledge gained from a pre-trained model on a large dataset for a specific task and applying it to a different, but related, task. In our case, we used the MobileNet model, which was pre-trained on a vast dataset, as a starting point. This allowed our model to inherit knowledge about general image features, enabling it to focus on the intricacies of brain tumor classification.

Fine-tuning played a crucial role in tailoring the pre-trained MobileNet model to our specific medical imaging task. We incorporated dense and dropout layers while adjusting various hyperparameters to optimize the model's performance. The addition of these layers facilitated better feature extraction and prevented overfitting, contributing to the remarkable validation accuracy of 97.24%, test accuracy of 97.86%, and an F1 score of 97.86%.

To provide transparency and insights into our model's decision-making process, we adopted GradCam for visual explanations. This not only aids medical professionals in understanding the model's predictions but also accelerates medical examinations and treatments, making them more efficient and cost-effective.

On the other hand, while our project has shown promising results, certain drawbacks warrant consideration. One significant limitation is the size of the dataset used for training the model. The availability of a relatively small dataset can hinder the model's ability to generalize effectively to diverse and unseen cases. To address the issue of a small dataset, a potential solution involves acquiring and incorporating a more extensive and diverse set of MRI images for training. Collaborating with multiple medical institutions to aggregate data or exploring the use of data augmentation techniques could help augment the dataset, providing the model with a richer understanding of the variations in brain tumor presentations.

Additionally, the current model may face challenges in precisely identifying the boundaries of tumors, potentially leading to false positives or negatives. Furthermore, the reliance on a pre-trained MobileNet model, while beneficial for leveraging general image features, may introduce biases or limitations in capturing subtle nuances unique to medical images. Developing a custom architecture tailored to the intricacies of medical imaging, perhaps through architecture search techniques, could lead to a more specialized and optimized model for brain tumor classification.

In the future, the future trajectory involves refining data preparation, adopting advanced visualization methods, and expanding the dataset. We plan to enhance our model by refining our data preparation techniques, employing advanced visualization methods, and expanding our dataset. By doing so, we aim to further increase the accuracy and robustness of our model, reinforcing its role as a valuable tool in the medical field for the accurate and prompt classification of brain tumors in MRI scans. Our ongoing efforts underscore the significance of artificial intelligence in advancing medical diagnostics and treatment processes.

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] Q. T. Ostrom, M. Price, C. Neff, G. Cioffi, K. A. Waite, C. Kruchko, and J. S. Barnholtz-Sloan, "CBTRUS Statistical Report: Primary Brain and Other Central Nervous System Tumors Diagnosed in the United States in 2016—2020," Neuro-Oncology, Oxford University Press US, vol. 25, no. Supplement_4, pp. iv1–iv99, 2023.

[2] J. Huang, H. Li, H. Yan, F.-X. Li, M. Tang, and D.-L. Lu, "The comparative burden of brain and central nervous system cancers from 1990 to 2019 between China and the United States and predicting the future burden," Frontiers in Public Health, Frontiers Media SA, vol. 10, p. 1018836, 2022.

[3] Q. T. Ostrom, N. Patil, G. Cioffi, K. Waite, C. Kruchko, and J. S. Barnholtz-Sloan, "CBTRUS statistical report: primary brain and other central nervous system tumors diagnosed in the United States in 2013–2017," Neuro-Oncology, Oxford University Press US, vol. 22, no. Supplement_1, pp. iv1–iv96, 2020.

[4] Q. T. Ostrom, M. Price, C. Neff, G. Cioffi, K. A. Waite, C. Kruchko, and J. S. Barnholtz-Sloan, "CBTRUS statistical report: primary brain and other central nervous system tumors diagnosed in the United States in 2015–2019," Neuro-Oncology, Oxford University Press US, vol. 24, no. Supplement_5, pp. v1–v95, 2022.

[5] G. Varoquaux and V. Cheplygina, "Machine learning for medical imaging: methodological failures and recommendations for the future," NPJ Digital Medicine, Nature Publishing Group UK London, vol. 5, no. 1, p. 48, 2022.

[6] M. Injadat, A. Moubayed, A. B. Nassif, and A. Shami, "Machine learning towards intelligent systems: applications, challenges, and opportunities," IEEE Transactions on Artificial Intelligence Review, Springer, vol. 54, pp. 3299–3348, 2021.

[7] X. Yu, J. Wang, Q.-Q. Hong, R. Teku, S.-H. Wang, and Y.-D. Zhang, "Transfer learning for medical image analyses: A survey," IEEE Transactions on Neural Networks and Learning Systems, Elsevier, vol. 489, pp. 230–254, 2022.

[8] E. Radiya-Dixit and X. Wang, "How fine can fine-tuning be? Learning efficient language models," in Proceedings of the International Conference on Artificial Intelligence and Statistics, PMLR, pp. 2435–2443, 2020.

[9] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, "A survey of convolutional neural networks: analysis, applications, and prospects," IEEE Trans. Neural Netw. Learn. Syst., 2021.

[10] W. Ayadi, W. Elhamzi, and M. Atri, "A new deep CNN for brain tumor classification," in 2020 20th International Conference on Sciences and Techniques of Automatic Control and Computer Engineering (STA), IEEE, pp. 266–270, 2020.

[11] A. Saleh, R. Sukaik, and S. S. Abu-Naser, "Brain tumor classification using deep learning," in 2020 International Conference on Assistive and Rehabilitation Technologies (iCareTech), IEEE, pp. 131–136, 2020.

[12] W. Ayadi, I. Charfi, W. Elhamzi, and M. Atri, "Brain tumor classification based on hybrid approach," The Visual Computer, vol. 38, no. 1, pp. 107-117, 2022.

[13] M. I. Sharif, M. A. Khan, M. Alhussein, K. Aurangzeb, and M. Raza, "A decision support system for multimodal brain tumor classification using deep learning," Complex & Intelligent Systems, pp. 1-14, 2021.

[14] S. Kumar and D. P. Mankame, "Optimization driven deep convolution neural network for brain tumor classification," Biocybernetics and Biomedical Engineering, vol. 40, no. 3, pp. 1190-1204, 2020.

[15] F. J. Díaz-Pernas, M. Martínez-Zarzuela, M. Antón-Rodríguez, and D. González-Ortega, "A deep learning approach for brain tumor classification and segmentation using a multiscale convolutional neural network," in Healthcare, vol. 9, no. 2, p. 153, 2021.

[16] H. T. Nguyen, H. H. Luong, P. T. Phan, H. H. D. Nguyen, D. Ly, D. M. Phan, and T. T. Do, "HS-UNET-ID: An approach for human skin classification integrating between UNET and improved dense convolutional network," International Journal of Imaging Systems and Technology, vol. 32, no. 6, pp. 1832-1845, 2022.

[17] H. H. Luong, N. H. Khang, N. Q. Le, D. M. Canh, P. S. Ha, et al., "A Proposed Approach for Monkeypox Classification," International Journal of Advanced Computer Science and Applications, vol. 14, no. 8, 2023

[18] H. T. Nguyen, H. H. Luong, T. H. N. Kien, N. P. L. Phan, T. D. Thuan, T. T. Tin, T. D. Nguyen, and T. C. Toai, "Brain Tumors Detection on MRI Images with K-means Clustering and Residual Networks," in International Conference on Computational Collective Intelligence, pp. 317-329, 2022.

[19] S. Rinesh, K. Maheswari, B. Arthi, P. Sherubha, A. Vijay, S. Sridhar, T. Rajendran, Y. A. Waji, et al., "Investigations on brain tumor classification using hybrid machine learning algorithms," Journal of Healthcare Engineering, vol. 2022, 2022.

[20] Z. N. K. Swati, Q. Zhao, M. Kabir, F. Ali, Z. Ali, S. Ahmed, and J. Lu, "Brain tumor classification for MR images using transfer learning and fine-tuning," Computerized Medical Imaging and Graphics, vol. 75, pp. 34-46, 2019.

[21] N. Abiwinanda, M. Hanif, S. T. Hesaputra, A. Handayani, and T. R. Mengko, "Brain tumor classification using convolutional neural network," in World Congress on Medical Physics and Biomedical Engineering 2018: June 3-8, 2018, Prague, Czech Republic (Vol. 1), pp. 183-189, 2019.

[22] H. H. Luong, L. T. T. Le, H. T. Nguyen, V. Q. Hua, K. V. Nguyen, T. N. P. Bach, T. N. A. Nguyen, and H. T. Q. Nguyen, "Transfer learning with fine-tuning on MobileNet and Grad-CAM for bones abnormalities diagnosis," in Computational Intelligence in Security for Information Systems Conference, pp. 171-179, 2022.

[23] H. H. Luong, N. T. L. Phan, T. C. Dinh, T. M. Dang, T. T. Duong, T. D. Nguyen, and H. T. Nguyen, "Fine-Tuning MobileNet for Breast Cancer Diagnosis," in Inventive Computation and Information Technologies: Proceedings of ICICIT 2022, pp. 841-856, 2023.

[24] H. H. Luong, P. T. Vo, H. C. Phan, N. L. D. Tran, H. Q. Le, and H. T. Nguyen, "Fine-Tuning VGG16 for Alzheimer's Disease Diagnosis," in Conference on Complex, Intelligent, and Software Intensive Systems, pp. 68-79, 2023.

# DDoS Classification using Combined Techniques

Mohd Azahari Mohd Yusof[1], Noor Zuraidin Mohd Safar[2], Zubaile Abdullah[3], Firkhan Ali Hamid Ali[4],
Khairul Amin Mohamad Sukri[5], Muhamad Hanif Jofri[6], Juliana Mohamed[7], Abdul Halim Omar[8],
Ida Aryanie Bahrudin[9], Mohd Hatta Mohamed Ali @ Md Hani[10]

College of Computing, Informatics and Mathematics,
Universiti Teknologi MARA (UiTM) Cawangan Melaka Kampus Jasin, Malaysia[1]
Faculty of Computer Science & Information Technology (FSKTM),
Universiti Tun Hussein Onn Malaysia (UTHM), Parit Raja, Batu Pahat, Johor, Malaysia[2, 3, 4, 5]
ICT as Enabler (iCAN) Focus Group, Department of Information Technology, Center for Diploma Studies,
Universiti Tun Hussein Onn Malaysia (UTHM), Pagoh Higher Education Hub, 84600 Pagoh, Johor, Malaysia[6, 7, 8, 9, 10]

*Abstract*—Now-a-days, the attacker's favourite is to disrupt a network system. An attacker has the capability to generate various types of DDoS attacks simultaneously, including the Smurf attack, ICMP flood, UDP flood, and TCP SYN flood. This DDoS issue encouraged the design of a classification technique against DDoS attacks that enter a computer network environment. The technique is called Packet Threshold Algorithm (PTA) and is combined with several machine learning to classify incoming packets that have been captured and recorded. Apart from that, the combination of techniques can differentiate between normal packets and DDoS attacks. The performance of all techniques in the research achieved high detection accuracy while mitigating the issue of a high false positive rate. The four techniques focused in this research are PTA-SVM, PTA-NB, PTA-LR and PTA-KNN. Based on the results of detection accuracy and false positive rate for all the techniques involved, it proves the PTA-KNN technique is a more effective technique in the context of detection of incoming packets whether DDoS attacks or normal packets.

*Keywords*—*DDoS; machine learning; accuracy; false positive rate*

## I. INTRODUCTION

The world now desperately needs an Internet to share resources with other users no matter where they are. It provides many facilities for users to perform daily activities including online games, social media and information search related to teaching and learning. Internet is available 24 hours a day to all users. However, the Internet is often threatened by several network attacks from attackers around the world and this includes DDoS attacks as said by study [1].

When a DDoS attack is launched by an attacker, the computer network or system is inaccessible at that time, even for users who have registered in the system. Typically, attackers apply botnets to perform DDoS attacks to get attacks with incredible speed. It can weaken the target server to serve all requests at that time. According to research [2], DDoS attacks can be categorized into three groups. These categories are volume-based attacks, followed by protocol attacks, and application layer attacks. Volume-based attacks are a category that involves attacks aiming to overwhelm network resources by flooding communication channels with a high volume of traffic. Volume-based attacks often utilize botnets, which are networks of compromised computers controlled by the attacker

[3]. By leveraging thousands or millions of infected devices, the attacker generates a massive amount of network traffic, leading to system failures in the targeted infrastructure. The category of protocol attacks focuses on attacking network protocol layers. DDoS protocol attacks often exploit vulnerabilities within the communication protocols used in network infrastructure, such as TCP/IP [4]. Attackers may employ techniques like SYN floods, where they send an overwhelming number of SYN requests to the target server, causing an overload of requests and hindering the server's ability to serve legitimate users. Meanwhile, application layer attacks refer to targeting specific applications or services running on top of the network infrastructure [5]. Application layer DDoS attacks focus on exploiting vulnerabilities within the application's logic or resources it relies on. Attackers can generate various types of DDoS attacks from anywhere. An example of such an attack is the HTTP flood, where attackers overwhelm a web server by sending an abnormally high volume of HTTP requests. This flood of requests leads to a strain on server resources, causing a degradation in performance or even a complete service failure.

There are several types of DDoS attacks that can be generated by attackers from anywhere. These attacks encompass ICMP flood, UDP flood, Ping of Death, Slowloris, Zero-day attack, Smurf, and TCP SYN flood [6]. In order to protect against DDoS attacks, a robust and effective detection strategy is crucial.

The research presents several significant contributions in the following manner:

- In this research, a DDoS attack classification algorithm called the Packet Threshold Algorithm (PTA) was developed to accurately distinguish incoming packets as either normal or malicious. It specifically targets TCP SYN flood, Smurf, UDP flood, Ping of Death, or normal packets. The PTA utilizes a packet threshold mechanism to differentiate and classify incoming traffic.

- To enhance detection capabilities and address the issue of false positives, the PTA was combined with various machine learning techniques, including Support Vector Machine (SVM), K-Nearest Neighbor (KNN), Naïve Bayes (NB), and Logistic Regression (LR). The

objective of this combination was to mitigate the problem of incorrectly classifying DDoS packets as normal packets or vice versa, as experienced in previous techniques.

- In addition to integrating machine learning algorithms with the PTA to enhance the overall performance and accuracy of the detection system, this research also conducted a comprehensive evaluation of the effectiveness of each technique. The results were presented to identify the most efficient approach for detecting malicious packets within a network environment. Furthermore, this study explored potential enhancements and optimizations to further advance the state-of-the-art in DDoS attack detection.

Having a reliable and precise detection strategy is of utmost importance in safeguarding against DDoS attacks. The combined approach of the PTA and machine learning algorithms significantly enhances the system's capability to accurately differentiate and classify incoming packets. By reducing false positives, this strategy provides a more effective defense against DDoS attacks, ensuring the integrity and availability of network resources [7].

The DDoS detection problem is enhanced using machine learning models such as SVM, KNN, Naïve Bayes, and Logistic Regression, which are well-suited for handling classification jobs. Naïve Bayes is strong at probabilistic classification, SVM is good at separating data points, KNN is good at pattern recognition, and Logistic Regression is good for binary classification. By adjusting to a variety of packet behaviors, these models help distinguish between malicious and legitimate packets with accuracy.

There are, nevertheless, certain restrictions. Large datasets may be a problem for SVM and KNN, affecting computing efficiency. The independence between features assumed by Naïve Bayes may not hold true for complex packet dynamics. Non-linear correlations between features may be difficult for logistic regression to handle, which could reduce its accuracy for complex packet classifications. When selecting the best model for DDoS detection, these limitations must be considered.

This paper is divided into several sections. Related work is presented in Section II. Next, in Section III, the methodology is presented, and the evaluation of techniques is described in Section IV, followed by results and discussion in Section V. The final section, Section VI, provides a brief summary of this paper.

## II. RELATED WORK

Despite the substantial research efforts dedicated to countering DDoS attacks, the challenge of mitigating them endures. Researchers have introduced various techniques in their attempts to combat the actions of DDoS attackers. Table I provides a summary of the methods proposed by these researchers to address such attacks.

Starting with the first study conducted by study [8], this research addresses the pressing issue of distributed denial-of-service (DDoS) attacks within the context of 5G networks. It emphasizes the predominant focus of previous studies on radio access networks (RAN) and voice service networks, often overlooking the vulnerabilities inherent in core networks (CN). These core network components, including the Access and Mobility Management Function (AMF), Session Management Function (SMF), and User Plane Function (UPF), are pivotal in providing expansive 5G coverage but are susceptible to DDoS attacks. The study introduces a methodology and a threat detection system tailored to counter signalling DDoS attacks specifically targeting 5G standalone CNs. By leveraging fundamental machine learning classifiers and preprocessing techniques such as entropy-based analysis (EBA) and statistics-based analysis (SBA), the research demonstrates the effectiveness of proactive defense strategies against these attacks. Notably, the results underscore the RF classifier as the top performer, achieving an impressive average accuracy of 98.7%.

The second study, led by [9], underscores the critical role of the internet as a fundamental communication tool in contemporary society. In tandem with the internet's indispensability, the frequency and severity of cyber-attacks have escalated, with DDoS attacks ranking among the top five most impactful and costly cyber threats. DDoS attacks disrupt legitimate users' access to network resources, necessitating the development of swift and accurate detection methods to mitigate their considerable damage. The study adopts machine learning classification algorithms, including LR, DT, RF, Ada Boost, Gradient Boost, KNN, and NB to detect DDoS attacks using the CICDDoS2019 dataset, encompassing eleven distinct DDoS attack types characterized by 87 features. The research evaluates classifier performance through various metrics, revealing that AdaBoost and Gradient Boost excel in classification, while LR, KNN, and NB also exhibit strong performance. However, DT and RF classifiers demonstrate less effective classification results.

The third study, conducted by [10], addresses the ongoing challenge of effectively managing DDoS attacks, which pose a significant threat to network security by inundating target networks with malicious traffic from multiple sources. Despite the availability of various conventional methods for detecting DDoS attacks, rapidly identifying these threats using feature selection algorithms remains a formidable task. In this study, a hybrid approach is introduced, incorporating feature selection techniques such as chi-square, Extra Tree, and ANOVA, in conjunction with four machine learning classifiers: FR, DT, KNN, and XGBoost. The primary goal is to enable early detection of DDoS attacks on IoT devices. To validate the proposed methodology, the research employs the CICDDoS2019 dataset, which encompasses a wide range of DDoS attacks, and conducts assessments in a cloud-based environment (Google Colab). The experimental results demonstrate the superior performance of the hybrid methodology, achieving an impressive 82.5% reduction in features and attaining 98.34% accuracy with ANOVA for XGBoost, thereby facilitating the early identification of DDoS attacks on IoT devices.

The fourth study, conducted by study [11], pioneers a comprehensive approach to address pressing security concerns in IoT networks, with a specific focus on the persistent threat

posed by DDoS attacks. Their innovative solution involves the integration of SDN with IoT to reinforce security measures and access control. Despite this integration, DDoS attacks continue to pose a formidable challenge. To tackle this issue head-on, the study introduces an advanced machine learning-based security framework. They meticulously craft a controlled testing environment for simulating DDoS attacks, capturing network logs, preprocessing them into a structured dataset, and employing a trio of robust algorithms, namely NB, DT, and SVM for network packet classification. Remarkably, their framework attains impressive accuracy rates, achieving 97.4% for NB, 96.1% for SVM, and an outstanding 98.1% for DT, unequivocally showcasing its effectiveness in mitigating DDoS threats while optimizing resource utilization and proficiently managing network traffic. This pioneering approach holds substantial promise for elevating the security posture of IoT networks.

TABLE I.     PAST STUDY MACHINE LEARNING TECHNIQUE

| Title of Paper / Year of Published | Machine Learning DDoS Detection Techniques | | | |
|---|---|---|---|---|
| | *SVM* | *NB* | *LR* | *KNN* |
| Machine Learning Based Signalling DDoS Detection System for 5G Stand Alone Core Network (2022) | ✓ | ✓ | ✗ | ✗ |
| Detection of DDoS Attacks Using Machine Learning Classification Algorithms (2022) | ✗ | ✓ | ✓ | ✓ |
| Analysis of Machine Learning Classifiers for Early Detection of DDoS Attacks on IoT Devices (2022) | ✗ | ✗ | ✗ | ✓ |
| Towards a Machine Learning-Based Framework for DDoS Attack Detection in Software-Defined IoT (SD-IoT) Networks (2023) | ✓ | ✓ | ✗ | ✗ |
| Detection of DDoS Attack in IoT Traffic using Ensemble Machine Learning Techniques (2023) | ✓ | ✓ | ✓ | ✗ |

In the final study conducted by study [12], the focus is on investigating DDoS attacks within the context of the IoT. The research utilizes machine learning classifiers, including both bagging, and boosting techniques, to categorize attack traffic, making use of the CICDDoS2019 dataset designed to simulate DDoS attacks on the UDP and TCP protocols commonly employed in IoT networks. To tackle data imbalance, the study employs an ensemble sampling approach that combines random under-sampling and ADASYN oversampling. Feature selection is carried out using two methods: the Pearson correlation coefficient and the Extra Tree classifier. The results reveal that RF performs the best with minimal training and prediction time, and Extra Trees for feature selection outperforms the Pearson correlation coefficient method in terms of overall time efficiency for most classifiers. However, it's noteworthy that when using the Pearson correlation coefficient for feature selection, RF remains the optimal choice for attack detection.

After conducting an extensive analysis of prior research in the field of DDoS detection using machine learning methods, it becomes evident that there is a pressing need to improve the process of feature selection in the datasets utilized. It is of paramount importance to minimize the occurrence of false

positives in order to achieve a heightened level of detection precision. This revelation underscores the critical importance of carefully selecting relevant and efficient features for incorporation into DDoS detection and classification methodologies. By enhancing feature selection techniques, the potential for generating false positive alerts can be significantly reduced, resulting in outcomes that are more reliable and precise.

## III. PROPOSED METHODOLOGY

This section introduces the research methodology, which is organized into four phases as illustrated in Fig. 1, and it outlines various research activities.



Fig. 1.    Methodology of proposed DDoS detection.

### A. Dataset Preparation

A dataset containing several types of DDoS attacks and normal packets is provided in the first phase, as shown in Fig. 2. The dataset is relevant to research activities as it records multiple incoming packets, which are the primary focus.



Fig. 2.    Sample of DDoS dataset.

It includes various features such as source address, destination address, packet type, packet size, and packet class. For instance, the source address refers to the IP address of the sender generating the packet or traffic, while the destination address represents the IP address that receives the packets or traffic.

### B. Data Preprocessing

The second research phase is data preprocessing. This phase is crucial in research work as it requires expertise to transform the data into a comprehensible format. Two activities were conducted in this phase: data cleaning and data reduction. Data cleaning is indeed the first activity in the research process, as presented in Fig. 3. This method is called identification of missing values, which is utilized in the research. It indicates that if there is a missing value, the output

will show a value 1, 2, and so on. This means that there are missing values or empty cells in the Src_Addrs, Pkt_ID, and From_Node columns in the dataset used.

```
In [1]:  import pandas as pd
         ddos_df = pd.read_csv('Desktop/original_ddos_dataset.csv')
         ddos_df.isnull().sum()

Out[1]:  Src_Addr       2
         Dst_Addr       0
         Pkt_ID         1
         From_Node      1
         To_Node        0
         Packet_Type    0
         Pkt_Size       0
```

Fig. 3. Identification of missing values.

The second activity involves data reduction, reducing the number of data samples by identifying and eliminating duplicate rows in the dataset, as presented in Fig. 4.

```
In [1]:  import pandas as pd
         ddos_df = pd.read_csv('Desktop/original_ddos_dataset.csv')
         ddos_df.head(10).duplicated()

Out[1]:  0    False
         1    False
         2    False
         3     True
         4    False
         5    False
         6     True
         7    False
         8    False
         9    False
```

Fig. 4. Identification of duplicate data.

In this case, data duplication occurs in rows 3 and 6, which need removal to generate high-quality data and facilitate analysis. Both activities assist in obtaining complete, consistent, and high-quality data within the dataset.

### C. Data Splitting

In the third phase of the research, known as data splitting, further investigation proceeds. The dataset, consisting of a total of 240,000 samples, is partitioned into two distinct sets: the training set and the testing set, as outlined in Table II.

The training set plays a crucial role in assessing the effectiveness of machine learning methods by utilizing data samples from the dataset. On the other hand, the testing set is employed to evaluate these methods. The train and test functions were formed to separate these two sets of data. The dataset was divided according to the data distribution outlined in Table II. For example, the data separation for 80: 20 ratios allocates 80% for the training set and the remaining 20% for the testing set.

TABLE II. DATA SPLITTING (TRAINING:TESTING)

| No. | Data Splitting Training:Testing | No. of Samples | |
|---|---|---|---|
| | | Training | Testing |
| 1 | 50:50 | 120,000 | 120,000 |
| 2 | 60:40 | 144,000 | 96,000 |
| 3 | 70:30 | 168,000 | 72,000 |
| 4 | 80:20 | 192,000 | 48,000 |

### D. Packet Classification

Quality data has been selected, and this research continues with the final phase, which is packet classification. In this phase, a technique called Packet Threshold Algorithm (PTA) has been proposed. This PTA is able to identify incoming packets whether normal packets or DDoS attacks. PTA is combined with several machine learning techniques, SVM, KNN, NB and LR. In the research, the functioning of this PTA was analyzed, as shown in Fig. 5. First, the PTA will check incoming packets based on a predefined packet threshold, which involves packet size and packet type received by the server. If the received packet is TCP or UDP or ICMP and a size of less than 60 bytes per second, PTA will issue the incoming packet category is normal packet. If the server receives TCP packets larger than 60 bytes per second, PTA will issue the incoming packet category is TCP SYN flood. Meanwhile, if the server receives a packet size exceeding 60 bytes per second and carries UDP packets, the PTA will issue the incoming packet category is UDP flood. If the type of packet received by the server is an ICMP packet and the size exceeds 65,535 bytes per second, PTA will issue the incoming packet category is Ping of Death. Meanwhile, if the ICMP packet size is less than 65,535 bytes per second but exceeds 60 bytes per second, PTA will issue the incoming packet category is a Smurf attack. The PTA will act to drop all packets received by the server, for which the packet size exceeds 60 bytes per second and the PTA allows packet sizes less than 60 bytes per second to enter the network environment. Finally, PTA is combined with machine learning by involving several phases or activities including features selection, data splitting, construction and evaluation of the techniques involved.

Here is a summary of how PTA determines the category of incoming packets. Firstly, PTA utilizes a predefined packet threshold to evaluate incoming packets. Secondly, PTA examines the packet type and size to determine their respective categories, as described above. Finally, based on the determined category, PTA performs specific actions on the packet: dropping all packets received by the server that exceed 60 bytes per second and allowing packets with sizes less than 60 bytes per second to enter the network. By employing this approach, PTA can accurately classify incoming packets as normal or belonging to various types of DDoS attacks.

```
Step 1: Start
Step 2: Check incoming packets
        if (packet size < 60) and (packet type = TCP) or (packet type = UDP) or (packet type = ICMP) then
            packet class = Normal
        if (packet size ≥ 60) and (packet type = TCP) then
            packet class = TCP SYN flood
        if (packet size ≥ 60) and (packet type = UDP) then
            packet class = UDP flood
        if (packet size ≥ 65,535) and (packet type = ICMP) then
            packet class = Ping of Death
        if (packet size ≥ 60 and packet size < 65,535) and (packet type = ICMP) then
            packet class = Smurf
Step 3: Features selection (x for input and y for target)
Step 4: Split data into training and testing set
Step 5: Build technique
Step 6: Train and test the technique
Step 7: Evaluation (TP, FP, TN and FN)
Step 8: End
```

Fig. 5. Packet Threshold Algorithm (PTA).

## IV. EVALUATION OF TECHNIQUES

During the evaluation phase, detection accuracy and false positive rate are employed as metrics to analyze the precise number of packets detected by PTA and the occurrence of erroneous detections. This encompasses cases where normal packets are wrongly identified as DDoS attacks and instances where DDoS attacks are mistakenly classified as normal packets. The calculation of detection accuracy and false positive rate follows a widely accepted standard formula.

$$Accuracy = \frac{TP+TN}{TP+TN+FN+FP} \times 100 \qquad (1)$$

$$FPR = \frac{FP}{FP+TN} \times 100 \qquad (2)$$

The formula explanation above can be summarized as follows:

- True Positive (TP): Instances where the model correctly detected DDoS attacks when they occurred.

- False Negative (FN): Instances where the model failed to detect DDoS attacks when they were happening.

- False Positive (FP): Instances where the model incorrectly flagged normal traffic as DDoS attacks.

- True Negative (TN): Instances where the model correctly identified normal traffic as not being DDoS attacks.

These four evaluations can be illustrated using the confusion matrix in Table III. The confusion matrix is a crucial tool in machine learning, providing a detailed breakdown of a model's performance by categorizing predictions into TP, FN, FP, and TN. This breakdown helps assess both accuracy and the model's ability to identify positive and negative cases accurately. It is a fundamental instrument for improving classification model effectiveness in various domains, including DDoS attack detection.

TABLE III. CONFUSION MATRIX

| | | Predicted DDoS | |
|---|---|---|---|
| | | DDoS | Normal |
| **Actual DDoS** | DDoS | TP | FN |
| | Normal | FP | TN |

## V. RESULT AND DISCUSSION

In this section, the experimental results for the various techniques employed are presented. Starting with an evaluation of the effectiveness of the proposed method for DDoS attack detection, followed by a comparative analysis with previously utilized techniques.

### A. Performance Comparison of PTA with Machine Learning Techniques

This section presents the performance results for four combinations of PTA techniques with machine learning based on data splitting between training and testing sets, as shown in Table III. Upon analyzing the performance of each technique using a 50:50 data splitting, it becomes evident that the PTA-

KNN technique attains the highest detection accuracy of 99.86%. It is closely followed by the PTA-SVM technique, which also achieves a detection accuracy of 99.86%. The PTA-LR technique achieves a detection accuracy of 99.12%, whereas the PTA-NB technique reaches a detection accuracy of 98.70%.

Shifting focus to the 60:40 data splitting, the PTA-KNN technique once again emerges as the frontrunner, achieving the highest detection accuracy of 99.86%. Remarkably, the PTA-KNN technique surpasses the detection accuracies achieved by the PTA-SVM, PTA-LR, and PTA-NB techniques, which are 99.66%, 99.17%, and 98.72% respectively. For the 70:30 data splitting, the PTA-KNN technique continues to outperform the other techniques with a detection accuracy of 99.84%. The PTA-SVM, PTA-LR, and PTA-NB techniques achieve respective detection accuracies of 99.65%, 99.16%, and 98.69%. Table IV shows the performance comparison of PTA with machine learning techniques. When considering the 80:20 data splitting, the PTA-KNN technique showcases an impressive detection accuracy of 99.83%, surpassing the PTA-SVM technique that achieves a detection accuracy of 99.63%. Furthermore, the PTA-LR technique demonstrates an impressive detection accuracy of 99.17%, whereas the PTA-NB technique achieves a slightly lower accuracy of 98.68%. Through meticulous examination, it can be deduced that the PTA-KNN technique showcases a remarkable efficacy in identifying incoming packets, regardless of their nature as DDoS attacks or normal packets. Observing the statistical outcomes presented in Fig. 6, which depict the effectiveness of the PTA-KNN technique in the research. This effectiveness stems from its utilization of packet type and size as key criteria. This conclusion is further supported by the exceptional detection accuracies achieved across various data splitting ratios: 99.86% for 50:50, 99.86% for 60:40, 99.84% for 70:30, and 99.83% for 80:20.

TABLE IV. PERFORMANCE COMPARISON OF PTA WITH MACHINE LEARNING TECHNIQUES

| Technique | Data Splitting (Training:Testing) | Detection Accuracy | False Positive Rate |
|---|---|---|---|
| **PTA-NB** | 50:50 | 98.70% | 1.10% |
| | 60:40 | 98.72% | 1.08% |
| | 70:30 | 98.69% | 1.10% |
| | 80:20 | 98.68% | 1.08% |
| **PTA-KNN** | 50:50 | 99.86% | 0.01% |
| | 60:40 | 99.86% | 0.01% |
| | 70:30 | 99.84% | 0.01% |
| | 80:20 | 99.83% | 0.02% |
| **PTA-SVM** | 50:50 | 99.66% | 0.01% |
| | 60:40 | 99.66% | 0.01% |
| | 70:30 | 99.65% | 0.01% |
| | 80:20 | 99.63% | 0.02% |
| **PTA-LR** | 50:50 | 99.12% | 0.26% |
| | 60:40 | 99.17% | 0.25% |
| | 70:30 | 99.16% | 0.27% |
| | 80:20 | 99.17% | 0.26% |

Fig. 6. Statistical outcomes of PTA-KNN technique.

Referring to Table V, it is noteworthy that the detection accuracies presented therein exceed the performance of alternative techniques, thus emphasizing the superiority of the PTA-KNN technique. The detection accuracy percentages for the PTA-KNN technique are determined based on the number of successfully detected incoming packets. For the 50:50 data splitting, 119,827 incoming packets were accurately detected, while 173 packets were misclassified. In the case of the 60:40 data splitting, the PTA-KNN technique successfully identified 95,863 incoming packets as valid, with 137 packets misclassified. Similarly, for the 70:30 data splitting, the technique detected 71,882 incoming packets correctly, but there were 118 misclassified packets. Lastly, for the 80:20 data splitting, the PTA-KNN technique successfully detected 47,920 incoming packets, with 80 packets being misclassified.

TABLE V. DETECTION RESULTS FOR DIFFERENT DATA SPLITTING RATIOS AND PACKET TYPES USING COMBINATION TECHNIQUES

| Technique | Data Splitting (Training:Testing) | No. of Incoming Packet Detected | | | | |
|---|---|---|---|---|---|---|
| | | Normal | Ping of Death | Smurf | TCP SYN Flood | UDP Flood |
| PTA-NB | 50:50 | 106,499 | 354 | 733 | 208 | 10,649 |
| | 60:40 | 85,233 | 287 | 591 | 166 | 8,495 |
| | 70:30 | 63,943 | 206 | 455 | 128 | 6,324 |
| | 80:20 | 42,643 | 133 | 299 | 83 | 4,209 |
| PTA-KNN | 50:50 | 107,165 | 354 | 952 | 229 | 11,127 |
| | 60:40 | 85,766 | 287 | 765 | 182 | 8,863 |
| | 70:30 | 64,361 | 206 | 581 | 141 | 6,593 |
| | 80:20 | 42,914 | 133 | 393 | 88 | 4,392 |
| PTA-SVM | 50:50 | 107,168 | 354 | 731 | 216 | 11,123 |
| | 60:40 | 85,769 | 287 | 589 | 173 | 8,859 |
| | 70:30 | 64,363 | 206 | 454 | 135 | 6,590 |
| | 80:20 | 42,916 | 133 | 298 | 85 | 4,389 |
| PTA-LR | 50:50 | 107,168 | 354 | 529 | 212 | 10,680 |
| | 60:40 | 85,769 | 287 | 418 | 170 | 8,557 |
| | 70:30 | 64,363 | 206 | 321 | 131 | 6,376 |
| | 80:20 | 42,916 | 133 | 220 | 84 | 4,249 |

## B. Performance Comparison between Proposed DDoS Detection Technique and Previous Techniques

This section presents a performance comparison between the proposed DDoS detection technique and existing methods, as displayed in Table VI. Within the provided table, which demonstrates performance comparisons in terms of detection accuracy for various techniques across different years of publication, it becomes evident that the highest and lowest accuracies vary significantly among the diverse techniques and algorithms employed. Notably, the proposed technique stands out with the highest overall accuracy of 99.86%, achieved using the KNN algorithm. However, it is essential to emphasize that the lowest accuracy values are somewhat dispersed. For instance, in the case of Park et al., the lowest accuracy values for LR and KNN are denoted as NA, indicating a lack of available data. In contrast, for other techniques, such as Gaur and Kumar, the lowest accuracy is attributed solely to the KNN algorithm, which attains an accuracy of 91.39%.

TABLE VI. PERFORMANCE COMPARISON BETWEEN PROPOSED DDoS DETECTION TECHNIQUE AND PREVIOUS TECHNIQUES

| Technique/Year of Published | Performance Comparison in Terms of Detection Accuracy | | | |
|---|---|---|---|---|
| | SVM | NB | LR | KNN |
| Park et al. (2022) | 98.76% | 87.61% | NA | NA |
| Dasari and Devarakonda (2022) | NA | 99.58% | 99.58% | 99.55% |
| Gaur and Kumar (2022) | NA | NA | NA | 91.39% |
| Bhayo et al. (2023) | 96.10% | 97.40% | NA | NA |
| Pandey and Mishra (2023) | 96.24% | 98.23% | 89.76% | NA |
| Proposed Technique (2023) | 99.66% | 98.72% | 99.17% | 99.86% |

Overall, the proposed technique appears to exhibit the highest accuracy across most algorithms, rendering it a promising approach for detection. Nevertheless, it is crucial to consider other factors, such as computational complexity and practical applicability, when selecting a technique for a specific problem.

## VI. CONCLUSION

The team has extensively researched the capabilities of the PTA technique in detecting both DDoS attacks and normal packets. This involves utilizing a predefined packet threshold that considers factors such as packet size and the specific packet types that attackers may generate. By integrating the PTA technique with diverse machine learning approaches, findings reveal that the PTA-KNN technique surpasses PTA-NB, PTA-SVM, and PTA-LR techniques in terms of detection accuracy and false positive rate percentage.

In the research, potential areas for future enhancement have also been identified based on findings. One possible direction for improvement involves exploring adaptive thresholding techniques that dynamically adjust the packet threshold based on network conditions and attack patterns. Additionally, investigating the integration of anomaly detection algorithms and deep learning models could enhance the PTA technique's ability to detect emerging and sophisticated DDoS attacks.

These avenues for future research aim to further enhance the effectiveness and resilience of the PTA technique in combatting evolving cyber threats.

REFERENCES

[1] X. Wang, Y. Li, H. J. Khasraghi, and C. Trumbach, "The Mediating Role of Security Anxiety in Internet Threat Avoidance Behavior," Computer Security, vol. 134, pp. 1–14, Nov. 2023, doi: 10.1016/j.cose.2023.103429.

[2] R. M. A. Haseeb-ur-rehman et al., "High-Speed Network DDoS Attack Detection: A Survey," Sensors, vol. 23, no. 15. Multidisciplinary Digital Publishing Institute (MDPI), pp. 1–25, Aug. 01, 2023. doi: 10.3390/s23156850.

[3] Y. Li and Q. Liu, "A Comprehensive Review Study of Cyber-Attacks and Cyber Security: Emerging Trends and Recent Developments," Energy Reports, vol. 7, pp. 8176–8186, Nov. 2021, doi: 10.1016/j.egyr.2021.08.126.

[4] N. Abosata, S. Al-Rubaye, G. Inalhan, and C. Emmanouilidis, "Internet of Things for System Integrity: A Comprehensive Survey on Security, Attacks and Countermeasures for Industrial Applications," Sensors, vol. 21, no. 11. MDPI AG, pp. 1–29, Jun. 01, 2021. doi: 10.3390/s21113654.

[5] N. Tripathi and N. Hubballi, "Application Layer Denial-of-Service Attacks and Defense Mechanisms: A Survey," ACM Computing Surveys, vol. 54, no. 4. Association for Computing Machinery, pp. 1–33, Jul. 01, 2021. doi: 10.1145/3448291.

[6] L. Zhou, Y. Zhu, Y. Xiang, and T. Zong, "A Novel Feature-Based Framework Enabling Multi-Type DDoS Attacks Detection," World Wide Web, vol. 26, no. 1, pp. 163–185, Jan. 2023, doi: 10.1007/s11280-022-01040-3.

[7] F. M. Salem, H. Youssef, I. Ali, and A. Haggag, "A Variable-Trust Threshold-Based Approach for DDoS Attack Mitigation in Software Defined Networks," PLoS One, vol. 17, no. 8, pp. 1–19, Aug. 2022, doi: 10.1371/journal.pone.0273681.

[8] S. Park, B. Cho, D. Kim, and I. You, "Machine Learning Based Signaling DDoS Detection System for 5G Stand Alone Core Network," Applied Sciences (Switzerland), vol. 12, no. 23, pp. 1–27, Dec. 2022, doi: 10.3390/app122312456.

[9] K. B. Dasari and N. Devarakonda, "Detection of DDoS Attacks Using Machine Learning Classification Algorithms," International Journal of Computer Network and Information Security, vol. 14, no. 6, pp. 89–97, Dec. 2022, doi: 10.5815/ijcnis.2022.06.07.

[10] V. Gaur and R. Kumar, "Analysis of Machine Learning Classifiers for Early Detection of DDoS Attacks on IoT Devices," Arabian Journal for Science and Engineering, vol. 47, no. 2, pp. 1353–1374, Feb. 2022, doi: 10.1007/s13369-021-05947-3.

[11] J. Bhayo, S. A. Shah, S. Hameed, A. Ahmed, J. Nasir, and D. Draheim, "Towards a Machine Learning-Based Framework for DDoS Attack Detection in Software-Defined IoT (SD-IoT) Networks," Engineering Applications of Artificial Intelligence, vol. 123, no. 1, pp. 1–17, Aug. 2023, doi: 10.1016/j.engappai.2023.106432.

[12] N. Pandey and P. K. Mishra, "Detection of DDoS Attack in IoT Traffic using Ensemble Machine Learning Techniques," Networks and Heterogeneous Media, vol. 18, no. 4, pp. 1393–1408, 2023, doi: 10.3934/nhm.2023061.

# Association Model of Temperature and Cattle Weight Influencing the Weight Loss of Cattle Due to Stress During Transportation

Jajam Haerul Jaman[1], Agus Buono[2], Dewi Apri Astuti[3], Sony Hartono Wijaya[4], Burhanuddin[5], Jajam Haerul Jaman[6]*

Fakultas Ilmu Komputer, Universitas Singaperbangsa Karawang, Karawang, Indonesia[1]
Deparemen Ilmu Komputer, IPB University, Bogor, Indonesia[2, 4, 6]
Deparemen Ilmu Nutrisi dan Ilmu Hewan, IPB University, Bogor, Indonesia[3]
Deparemen Agribusiness, IPB University, Bogor, Indonesia[5]

*Abstract*—This study aimed to enhance animal welfare in the context of modern agriculture. The Association Rule analysis method using FP-Growth and Apriori algorithms was employed to identify patterns and factors influencing animal welfare, particularly in the context of live cattle weight loss (shrink) due to stress during transportation. Data obtained from several farms and clinical tests were used to develop insights into the relationship between farming practices, data science, and animal welfare. The research stages included data preprocessing, initial analysis, modeling, evaluation, and interpretation of results, recommendations and implications, and conclusions. The research results indicate that the use of FP-Growth and Apriori algorithms uncovered hidden patterns in the data, resulting in four association rules from FP-Growth and five rules from Apriori. These rules aid in designing recommendations to enhance animal welfare, improve agricultural efficiency, and support sustainability of the cattle sector. Our findings have significant implications in the context of animal welfare and sustainable farm management.

*Keywords—Association rule; animal welfare; cattle management; animal product quality; modern agriculture; recommendations; sustainability*

## I. INTRODUCTION

The issue of tracking and handling pressure in cattle, especially farm animals, for the duration of transportation, has resulted in extensive difficulties for cattle enterprises. The transportation of stay cattle often entails long distances from production facilities to consumption centers, resulting in diverse, demanding situations that need to be addressed. In this context, pressure on farm animals through transportation has been proven to result in sizable weight reduction, which, in turn, impacts the first-rate of meat and promotes charges [1], [2]. Consequently, efforts to reduce and control stress in cattle throughout transportation are essential [3], [4]. One of the factors that can influence the stress levels and weight loss in cattle during transportation is the body temperature condition. This study aims to explore the relationship between body temperature and the body weight of cattle in affecting weight loss due to stress during transportation. The method involves measuring the body temperature of cattle before transportation. Additionally, the body weight of cattle is also measured before and after transportation to calculate the weight difference caused by stress during transportation [5].

Researchers and stakeholders in cattle enterprises have long sought to address this issue. Numerous techniques have been developed, ranging from gazing at animal conduct to measuring physiological and biochemical parameters [6]. However, an essential question arises: are the current techniques adequate and reliable? [7]. In this research, efforts are made to address this question using association techniques. This association method will examine the relationship between the body temperature of cattle and their body weight in influencing weight loss in cattle during transportation.

In this section, we explain the association rules and how they can be used to gain valuable insights into stress in cattle during transport will be provided [8], [9]. Related studies that have successfully employed association rules in various contexts, including traffic pattern comprehension and the detection of unusual events in videos, will be referred to by us [10], [11].

Furthermore, we can discuss the study findings associated with using the Apriori and FP-increase algorithms in determining affiliation guidelines in our dataset, which incorporates temperature statistics, initial weight, and the fee of weight loss of cattle at some point of transportation [12]. These two algorithms are also be compared to assess their effectiveness in generating informative association rules [13].

To provide a broader context, other related research that has employed data mining techniques to address various issues, such as library user behavior analysis and deformation response analysis in landslide hazards, will also be referred to by us [7]. The knowledge gained from these studies offers a broader perspective on the potential use of association rules in the context of stress management in live cattle during transportation in the cattle industry [14].

In conclusion, our findings will be summarized, and further recommendations and implications from this research will be provided to minimize the impact of stress on live cattle during transportation in the cattle industry [4], [15].

## II. RELATED WORK

The transportation of stray animals is a critical element in the cattle industry, especially when these animals are destined for slaughter or processing. Previous research has indicated

that animal transportation can negatively impact animal welfare [16]–[18]. Animal welfare was assessed via signs that included pressure tiers, weight reduction, and the physical condition of the animals. Several elements have been identified as likely to affect animal welfare throughout transportation. These elements included car density, environmental temperature, travel duration, and animal management. The author in [19] showed that excessive car density can increase the strain degree in animals. Weight loss at some point of transportation is a severe problem that could impact animal productivity and the quality of the resulting meat. Significant weight loss in animals during transport has been documented in several studies [20]–[22]. This may illustrate the stress and soreness experienced by animals during transportation. Governments and global agencies have introduced diverse policies and suggestions to address animal welfare issues at some point of transportation. For instance, the ecu Union has applied strict policies regarding the shipping of live animals [23][24]. However, the implementation and enforcement of these regulations can range across international locations and areas .

While numerous studies have been conducted in the realm of animal transportation, several unaddressed research areas remain. Some investigations may focus on specific animal types or geographic regions, whereas others maintain a broader scope. Furthermore, our understanding of the individual factors within animals that can influence their responses to transportation remains incomplete. In the context of this literature review, our research endeavors to bridge these knowledge gaps by concentrating on live cattle and the variables that affect their well-being and weight loss during transport. This study aimed to provide a comprehensive framework for research within the domain of live animal transportation, elucidating its relevance to our own research. Our findings are correlated with prior research outcomes to facilitate a more profound understanding of this issue and its potential to enhance animal welfare in the cattle transportation sector.

Investigations within the realm of animal transport, specifically concentrating on the interplay between temperature, body weight, and cattle stress, have attracted noteworthy interest. Previous research has extensively explored the nuanced relationships between these factors, offering valuable insights into the difficulties confronted by livestock during transit.

### A. *Korelasi Temperature and Stress Correlation*

Previous research has explored the impact of temperature on the stress levels of cattle during transportation. Findings indicate a significant correlation, where increased temperatures often contribute to heightened stress levels [25][5]. These studies utilized various methodologies, including real-time temperature monitoring and observations of stress behavior, to establish robust associations.

The results of these studies suggest that high temperatures can influence the comfort and well-being of cattle during transportation, leading to an increase in stress levels. Furthermore, the research also indicates that cattle experiencing high stress levels during transportation tend to

undergo more significant weight loss. Thus, there is a connection between the body temperature of cattle and body weight in influencing weight loss during transportation..

### B. *Body Weight Dynamics*

The influence of body weight on cattle stress during transportation has been a focal point in several studies. Researchers have examined how fluctuations in body weight, particularly weight loss, align with increased stress levels. Through comprehensive analysis, these studies aim to unveil patterns and associations contributing to a deeper understanding of stress dynamics during transit. This approach involves measuring the body weight of cattle before and after transportation, along with monitoring their body temperature conditions [6].

The research findings indicate that high temperatures can affect the comfort and well-being of cattle during transportation, leading to an increase in stress levels. Additionally, the study suggests that cattle experiencing high stress levels during transportation tend to undergo more significant weight loss. Therefore, there is a correlation between cattle body temperature and body weight in influencing weight loss during transportation.

## III. RESEARCH METHOD

The initial steps in this research process involved data collection, preliminary analysis, and data preprocessing. The data obtained, such as information on body temperature, initial weight of the cattle, and rate of weight loss during transportation, will serve as the primary foundation for advancing this research. Following this, the subsequent phases of our study will involve utilizing association rule algorithms to detect and comprehend the connections among these factors concerning the stress encountered by cattle during transportation. For a summary of the research process, please consult Fig. 1.



Fig. 1. Methodology.

### A. *Data Collection*

The primary dataset will encompass details concerning the cattle's initial body temperature, starting weight, and rate of weight loss throughout transportation, sourced from pertinent outlets within the cattle sector. Additionally, information regarding cattle behavior during transit will be documented, encompassing stress indicators like agitation, profuse sweating,

or any unusual conduct, as supplementary data for potential future use.

### B. Basic Analysis

A descriptive analysis will be conducted on the data to initially comprehend the distribution of pertinent variables. Visual aids like charts, histograms, or scatter plots will be employed to scrutinize data patterns. This was done to provide a better understanding of the data, identify the research significance, determine the most appropriate methods, develop hypotheses and theoretical frameworks, improve planning, avoid data redundancy, and enhance the validity and reliability of the research.

### C. Data Preprocessing

The collected data will be analyzed to identify and address missing or invalid data. Temperature, weight, and weight loss rate data were converted into appropriate formats for statistical analysis, including data cleaning, removal of duplicate data, data selection, data transformation, and data encoding if necessary.

### D. Implementation of Association Rule Algorithm

Association rule mining is a data-mining technique aimed at discovering implicit association rules from a large amount of transactional data. It is also known as association analysis, which establishes connections between various item combinations within a dataset [26]. This association model was applied to the dataset to identify association rules related to stress in cattle during transportation. This involves configuring parameters, such as support, confidence, and lift ratio. There are two association rule models that we will use: Apriori and the FP-Growth algorithm. The Apriori algorithm is an association rule algorithm commonly used to address issues in data analysis involving transactions or lists of items, such as in retail shopping, online shopping carts, product recommendations, or pattern detection in transactional data. The Apriori algorithm is used to discover association rules connecting items or attributes in transactional datasets. Similar to the Apriori algorithm, FP-Growth is a model algorithm for associations resulting from the development of the Apriori algorithm. The primary use cases for both of these algorithms usually involve solving Market Basket Analysis problems, providing the most suitable product recommendations, and analyzing customer purchase patterns. Additionally, FP-Growth has the capability to handle large datasets, analyze

transactional data, and provide recommendations in e-commerce. Here are the required notations to generate an association rule.

To measure the initial probabilities of X and Y, we need an equation called Support, which is as follows [27], [28]:

$$S_{supp}(X \to Y) = \frac{(|X \cup Y|)}{n} \qquad (1)$$

$S_{supp}$ = Support value

n = total number of transactions.

After obtaining support, the next step is to find the Confidence value, which is a proportion of transactions containing all items in both X and Y compared to transactions containing only items in X [29], The notation for this is as follows:

$$Confidence(X \to Y) = \frac{(X \cup Y)}{Support(X)} \qquad (2)$$

Lift(X → Y) is a measure used to assess the significance of the association rule X → Y. Support (X ∪ Y) reflects how often itemsets X and Y appear together in transactions, whereas Support (X) and Support (Y) are measures of how often each individual itemset appears in transactions.

### E. Evaluation and Interpretation of Results

The results from the use of both algorithms will be evaluated to measure their effectiveness in generating informative association rules. Significant association rules were interpreted to understand the relationship between temperature, initial weight, rate of weight loss, and stress in cattle during transportation. Based on the analysis results, this research will also summarize the findings and conclude that association rules can provide valuable insights into the measurement and management of stress in cattle during transportation in the cattle industry.

## IV. RESULT

### A. Prepare Dataset

Data collection was carried out through observation. Body temperature data before departure, initial weight, and rate of weight loss during the journey were the primary focus of this study. Fifty data records have been successfully collected from relevant sources in the cattle industry, and for more details, the dataset can be seen in Table I.

TABLE I.        DATASET OF PRE-SHIPPING LIVE CATTLE OBSERVATION RESULTS

| NO | x1 | x2 | x3 | x4 | SUHU AWAL | | | | | | | | x13 | x14 | Shr/kg | % |
|----|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|--------|---|
| | | | | | x5 | x6 | x7 | x8 | x9i | x10 | x11 | x12 | | | | y |
| 1 | PO | 231 | M | 4 | 36 | 37 | 36 | 36 | 36 | 36 | 35 | 38 | 36 | 205 | 26 | 11% |
| 2 | PO | 235 | M | 4 | 37 | 37 | 37 | 37 | 37 | 37 | 36 | 39 | 37 | 211 | 24 | 10% |
| 3 | PO | 246 | M | 5 | 37 | 37 | 37 | 37 | 37 | 37 | 36 | 39 | 37 | 235 | 11 | 4% |
| 4 | PO | 239 | M | 5 | 36 | 37 | 36 | 37 | 37 | 37 | 35 | 38 | 36 | 203 | 36 | 15% |
| 5 | PO | 231 | M | 4 | 36 | 37 | 37 | 37 | 37 | 37 | 36 | 37 | 37 | 188 | 43 | 19% |
| …. | …… | ….. | … | …. | …. | ….. | …. | …. | …. | …. | …. | …. | … | … | … | … |
| …. | …… | ….. | … | …. | …. | ….. | …. | …. | …. | …. | …. | …. | … | … | … | … |
| 49 | Bali | 151,06 | F | 4 | 37 | 37 | 40 | 40 | 38 | 39 | 40 | 39 | 39 | 144 | 7 | 5% |
| 50 | Bali | 143,81 | M | 4 | 38 | 39 | 40 | 40 | 38 | 39 | 40 | 39 | 39 | 136 | 7 | 5% |

The observation results yielded a dataset consisting of 50 data records with 14 independent variables (x) and 1 dependent variable (y). x1 represents the breed of cattle, with two attributes: Peranakan Ongole (PO) and Bali Cattle. x2 represents the initial weight of the cattle, measured either through weighing or estimation based on the chest circumference and body length. x3 indicates the sex of cattle identified through direct observation. Column D represents the age of the cattle. x4–x12 represent temperature measurements taken from various body parts (forehead, orbital, cheek, shoulder, back, thigh, leg, and rectal) using a thermogun for external temperature and a regular thermometer for rectal temperature. x13 represents the average body temperature of the cattle. x14 represents the final weight of the cattle upon reaching their destination. Shr/kg denotes the shrinkage in weight during transportation, measured in kilograms. Column y represents the percentage of weight shrinkage during transportation (Shrink). If observed, we only have 50 data records, a number that indeed appears limited, but obtaining them posed a very challenging task. We utilize this record count as part of algorithm testing, examining how well the algorithm can perform its task with such limited data. Additionally, the limitation in the number of records is turned into a strength in this research.

### B. Basic Analysis

Descriptive statistics for the 50 data records indicated that the measurement of live cattle weight before shipment (x2) resulted in a weight range of 125–261 kg, with an average of 195 kg. Age (x3) has a ranged from 3.5 to 4.77 years. The average temperature of the cattle (x13) before shipment range between 35.6°C and 37.4°C. The final weight (x14) has a ranges between 118 kg and 235 kg, with an average of 175.9 kg. The percentage of live cattle weight reduction during transportation (y) ranged from 0% to 27%. We also calculated the correlations for each variable. Table II and Fig. 2 show the results and visualization of the correlation between variables x and y, respectively.

TABLE II.    CORRELATION OF EACH VARIABLE WITH THE TARGET VARIABLE

|  | X2 | X4 | X5 | X6 | X7 | X8 | X9 | X10 | X11 | X12 | X13 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| CC | 0,33 | -0,22 | -0,50 | -0,28 | -0,35 | -0,43 | -0,41 | -0,29 | -0,40 | -0,22 | -0,49 |
| CP | 0,92 | 0,45 | -0,16 | -0,17 | -0,30 | -0,07 | 0,01 | -0,04 | 0,22 | -0,28 | -0,19 |

a. CC (Coeficient Corelation), CP (Coeficient Pearson)

From the initial analysis, it was found that there were some prunings on the variables that had a less significant impact on variable y. These include x1, x3 to x11, x13, and shr/kg. There are several reasons for this pruning, such as data type mismatches for x1 and x3. Additionally, for temperature measurements, discussions with stakeholders and experts led to an agreement to use only rectal temperature, whereas external body temperature measurements were deemed to have invalid calculations. The results of the initial analysis can be seen in Table III.



Fig. 2.    Illustration of the correlation between each variable x and the label y.

TABLE III.    DATASET RESULTING FROM VARIABLE PRUNING THAT IS READY FOR ANALYSIS

| No | Weight | Temp | Shrink |
|---|---|---|---|
| 1 | 231 | 38 | 11% |
| 2 | 235 | 39 | 10% |
| 3 | 246 | 39 | 4% |
| 4 | 239 | 38 | 15% |
| 5 | 231 | 37 | 19% |
| …. | ….. | …. | … |
| …. | ….. | …. | … |
| 49 | 151 | 39 | 5% |
| 50 | 144 | 39 | 5% |

Out of several variables that were pruned, only 2 predictor variables and 1 target variable were retained. This simplifies the subsequent analysis process, as there are not too many influential variables to be calculated.

### C. Data Preprocessing Process

In this process, we performed data cleaning. We also found that the data needed to be categorized to facilitate calculations, so we transformed the data into Crisp data. Rectal temperature was divided into three attributes: 1) Low Temperature: < 37.5°C (s_low); 2) Normal Temperature: 37.5°C – 39.5°C (s_normal); 3) High Temperature: > 39.5°C (s_high). The initial weight was divided into five attributes: 1) Very Thin Weight: 125 – 150 kg (w_very_thin); 2) Thin Weight: 151 – 175 kg (w_thin); 3) Medium Weight: 176 – 200 kg (w_medium); 4) Fat Weight: 201 – 225 kg (w_fat); 5) Very Fat Weight: 225 – 262 kg (w_very_fat). The percentage of shrinkage was divided into three attributes: 1) Slight Shrinkage Category: ≤ 9% (p_slight); 2) Moderate Shrinkage: 10% – 18% (p_moderate); 3) High Shrinkage: > 19% (p_high). The dataset resulting from the data transformation is shown in Table IV.

### D. Implementation of the Association Rule Algorithm

Descriptive statistics for the 50 data records indicated that the measurement of live cattle weight before shipment (x2) resulted in a weight range of 125–261 kg, with an average of 195 kg. Age (x3) has a ranged from 3.5 to 4.77 years. The

average temperature of the cattle (x13) before shipment range between 35.6°C and 37.4°C. The final weight (x14) has a ranges between 118 kg and 235 kg, with an average of 175.9 kg. The percentage of live cattle weight reduction during transportation (y) ranged from 0% to 27%. We also calculated the correlations for each variable. Table I and Fig. 1 show the results and visualization of the correlation between variables x and y, respectively.

TABLE IV.    THE RESULTS OF THE TRANSFORMATION BECOME CRISP DATA

| NO | WEIGHT | TEMP | SHRINK |
|---|---|---|---|
| 0 | w_very_thin | s_normal | p_slight |
| 1 | w_very_thin | s_normal | p_moderate |
| 2 | w_very_thin | s_normal | p_slight |
| 3 | w_very_thin | s_normal | p_slight |
| .... | ............................... | .................. | ................ |
| 18 | w_medium | s_normal | p_slight |
| 19 | w_thin | s_normal | p_moderate |
| .... | ............................... | .................. | ................ |
| 48 | w_fat | s_normal | p_slight |
| 49 | w_fat | s_normal | p_slight |

- Apriori Algorithm.

From the analysis conducted, relevant rules for the research were obtained. We set a minimum confidence level of 0.01 with a minimum threshold of 1.0. The results from using the Apriori algorithm produced 84 rules. Subsequently, we selected the most relevant rules from these, resulting in five rules that were most relevant to the actual conditions. The five rules that have been generated are listed in Table V.

TABLE V.    APRIORI ALGORITHM ASSOCIATION MODEL

| Antecedents | Consequents | Sup | Conf | Lift |
|---|---|---|---|---|
| weight_b_thin Temp_s_normal | Shrink_p_slight | 0,17 | 0,56 | 1,19 |
| weight_b_thin; Temp_s_High | Shrink_p_ slight | 0,07 | 1,00 | 2,14 |
| Temp_s_ normal; weight_b_medium | Shrink_p_moderate | 0,07 | 0,50 | 1,50 |
| Temp_s_low; weight_b_very_fat | Shrink_p_ moderate | 0,03 | 1,00 | 3,00 |
| Temp_s_ normal; weight_b_veri_thin | Shrink_p_slight | 0,13 | 0,80 | 1,71 |

The results obtained provide an overview that out of the five rules, they are supported by Support values ranging from 0.03 to 0.17, Confidence values from 0.50 to 1.00, and Lift Ratios from 1.19 to 3.00. Therefore, the average confidence values obtained served as the basis for the generated rules, indicating a relatively strong correlation.

- FP-Growth Algorithm.

Similarly to what was done previously with Apriori, we set a minimum confidence of 0.01 with a minimum threshold of 1.0 for the FP-Growth algorithm. The results of using the FP-Growth algorithm yielded 21 rules. We then conducted a

selection process for these rules, resulting in the identification of the four most relevant rules in line with the actual conditions. These four rules can be observed in Table VI.

TABLE VI.    FP-GROWTH ALGORITHM ASSOCIATION MODEL

| Antecedents | Consequents | Sup | Conf | Lift |
|---|---|---|---|---|
| Temp_s_low | Shrink_p_moderate | 0,03 | 1,00 | 3,00 |
| weight_w_very_thin'; 'Temp_s_low | Shrink_p_slight | 0,03 | 1,00 | 3,00 |
| Temp_s_high | Shrink_p_ slight | 0,07 | 1,00 | 2,14 |
| weight_b_thin; 'Temp_s_high | Shrink_p_ slight | 0,07 | 1,00 | 2,14 |

The results obtained indicate that out of the four rules, they are supported by Support values ranging from 0.03 to 0.07, Confidence values of 1.00, and Lift Ratio values ranging from 2.14 to 3.00. Therefore, the average Confidence value obtained serves as a strong basis for the generated rules.

The results of this research offer valuable insights into stress measurement and management in the transportation of live cattle within the cattle industry. The rules derived from the FP-Growth algorithm exhibited an exceptional confidence level of 1.00, signifying their robustness within the dataset and a high level of confidence in their applicability. In contrast, the Apriori algorithm yields rules with varying levels of Support and Confidence ranging from 0.50 to 1.00, highlighting the diversity in confidence levels among these rules. However, the overall self-assurance cost remains splendid, indicating a sturdy correlation with a number of popular policies. Furthermore, FP-growth generates rules with better increase Ratios than Apriori, indicating a higher correlation with a number of the rules it generates in the dataset. Tables VII and VIII show the rules given by the heatmap can further clarify the position of each rule.

TABLE VII.    RESULT HEATMAP APRIORI MODEL

| Weight / Temp | V Thin | Thin | Medium | Fat | Very Fat |
|---|---|---|---|---|---|
| Low | | | | | |
| Normal | | Slight S = 0,17 C=0,56 L = 1,19 | Moderate S = 0,07 C=0,50 L = 1,5 | Moderate S = 0,03 C=1,00 L=3,00 | Slight S = 0,13 C=0,56 L = 1,71 |
| High | | Slight S = 0,07 C=1,00 L = 2,14 | | | |

TABLE VIII.    RESULT HEATMAP FP-GROWTH MODEL

| Weight / Temp | Very Thin | Thin | Medium | Fat | Very Fat |
|---|---|---|---|---|---|
| Low | Moderate S = 0,03 C=1 L = 3 | Moderate S = 0,03 C=1 L = 3 | Moderate S = 0,03 C=1 L = 3 | Moderate S = 0,03 C=1 L = 3 | Moderate S = 0,03 C=1 L = 3 |
| Normal | | | | | |
| High | Slight S = 0,07 C=1 L = 2,14 | Slight S = 0,07 C=1 L = 2,14 | Slight S = 0,07 C=1 L = 2,14 | Slight S = 0,07 C=1 L = 2,14 | Slight S = 0,07 C=1 L = 2,14 |

TABLE IX.    RESULT HEATMAP APRIORI MODEL

| Confidence | 1 | 0,8 | 0,6 | 0,5 | 0,2 |
|---|---|---|---|---|---|
|  |  |  |  |  |  |

Four rules obtained from FP-Growth will automatically generate 10 rule. This is because rules 1 and 3 provide flexibility for any weight to be given its effect. Table IX shows the description of the heatmap by focusing on its confidence value.

## V.  CONCLUSIONS

The goal of this study was to determine the effect of temperature and initial weight on the pressure stages and weight reduction in cattle in the context of farm animal enterprises. The data collected and analyses conducted offer crucial insights into the factors affecting cattle well-being during transportation and their practical implications. Influence of Weight and Temperature on Cattle Stress. The evaluation results indicate that pre-transportation weight and temperature measurements significantly influence cattle pressure stages, that is, glaring within the association policies generated by each algorithm, which exhibit a strong correlation between weight and temperature and a lower frame weight (shrinkage) at the stop of transportation. This indicates that preliminary temperature and weight information can function as guidance for stakeholders to provide the maximum appropriate care of their farm animals, in the long run, lowering farm animal strain in the course of transportation. This information has significant practical implications for the cattle industry, allowing farmers to use temperature and weight data as vital indicators in managing cattle stress during transportation. Preventive measures, such as ensuring proper vehicle ventilation and safeguarding against animal welfare threats can assist in stress reduction and minimize detrimental weight loss in cattle.

This paper has certain limitations. Firstly, qualitative data was used to gather information on cattle behavior during transportation, and further analysis may be necessary to gain a deeper understanding of stress-related behavior. Secondly, the research focused solely on temperature as a contributing factor, disregarding other factors like humidity and population density within the transport vehicle, which might also influence cattle stress during transit. While this research serves as a critical foundation exploration of the factors influencing cattle stress during transportation, it is imperative to conduct subsequent studies with broader variables and larger sample sizes to offer a more comprehensive perspective and enhance stress management practices in the cattle industry.

## VI.  DISCUSSION

In this chapter, we will delve into the acquired research findings and provide additional context regarding their implications, along with offering recommendations for future research. Although this study has provided valuable insights, there are areas that warrant further investigation. Future research in this field should encompass several facets. Firstly, there is a need to collect more comprehensive behavioral data concerning animal behavior during transportation, as this study predominantly relied on qualitative data. Secondly, it is crucial to consider additional variables such as air humidity and

population density within the transport vehicle, as they may also exert an influence on cattle stress during journeys. Lastly, conducting advanced research with larger sample sizes and a broader spectrum of variables can provide a more holistic comprehension of the factors contributing to cattle stress during transportation, ultimately contributing to the enhancement of stress management practices within the cattle industry.

## REFERENCES

[1] K. Schwartzkopf-Genswein, J. Ahola, L. Edwards-Callaway, D. Hale, and J. Paterson, "Symposium Paper: Transportation issues affecting cattle well-being and considerations for the future," Prof. Anim. Sci., vol. 32, no. 6, pp. 707–716, 2016, doi: 10.15232/pas.2016-01517.

[2] G. Smith, C. State, and T. Grandin, "Effect of Transport on Meat Quality and Animal Welfare of Cattle , Pigs , Sheep , Horses , Deer , and Poultry," Anim. Welf., no. d, pp. 1–45, 2010.

[3] M. Alam et al., "Assessment of transport stress on cattle travelling a long distance (≈648 km), from Jessore (Indian border) to Chittagong, Bangladesh," Vet. Rec. Open, vol. 5, no. 1, pp. 1‑10, 2018, doi: 10.1136/vetreco-2017-000248.

[4] D. Ashenafi, E. Yidersal, E. Hussen, S. Tsegaye, and M. Desiye, "The Effect of long Distance Transportation Stress on Cattle: a Review," Biomed. J. Sci. Tech. Res., vol. 3, no. 3, pp. 3304–3308, 2018, doi: 10.26717/bjstr.2018.03.000908.

[5] G. E. Tresia, A. F. Trisiana, and B. Tiesnamurti, "Effects of Road Transportation on Some Physiological Stress Measures in Anpera and Boerka Goats," Bul. Peternak., vol. 47, no. 3, p. 142, 2023, doi: 10.21059/buletinpeternak.v47i3.83317.

[6] A. F. Trisiana, A. Destomo, and F. Mahmilia, "Pengangkutan Ternak : Proses, Kendala dan Pengaruhnya pada Ruminansia Kecil," Wartazoa, vol. 31, no. 1, pp. 43–53, 2021.

[7] L. Linwei, W. Yiping, H. Yepiao, L. Bo, M. Fasheng, and D. Ziqiang, "Optimized Apriori algorithm for deformation response analysis of landslide hazards," Comput. Geosci., vol. 170, no. October 2022, 2023, doi: 10.1016/j.cageo.2022.105261.

[8] M. C. Lucy, "Stress, strain, and pregnancy outcome in postpartum cows," Anim. Reprod., vol. 16, no. 3, pp. 455–464, 2019, doi: 10.21451/1984-3143-AR2019-0063.

[9] M. J. Prescott, Guidelines for the Human Transportation of Research Animals, vol. 28, no. 2. 2007. doi: 10.1007/s10764-007-9134-8.

[10] I. H. Sarker, "Machine Learning: Algorithms, Real-World Applications and Research Directions," SN Comput. Sci., vol. 2, no. 3, pp. 1–21, 2021, doi: 10.1007/s42979-021-00592-x.

[11] M. M. Taye, "Understanding of Machine Learning with Deep Learning: Architectures, Workflow, Applications and Future Directions," Computers, vol. 12, no. 5, 2023, doi: 10.3390/computers12050091.

[12] A. I. Idris et al., "Comparison of Apriori, Apriori-TID and FP-Growth Algorithms in Market Basket Analysis at Grocery Stores," IJICS (International J. Informatics Comput. Sci., vol. 6, no. 2, p. 107, 2022, doi: 10.30865/ijics.v6i2.4535.

[13] J. Wang and Z. Cheng, "FP-Growth based Regular Behaviors Auditing in Electric Management Information System," Procedia Comput. Sci., vol. 139, pp. 275–279, 2018, doi: 10.1016/j.procs.2018.10.268.

[14] Food and Agriculture Organization Of The United Nations, Livestock-Related Inventions During Emergencies, vol. 6, no. August. 2016.

[15] S. Larios-Cueto, R. Ramírez-Valverde, G. Aranda-Osorio, M. E. Ortega-Cerrilla, and J. C. García-Ortiz, "Indicadores de estrés en bovinos por el uso de prácticas de manejo en el embarque, transporte y desembarque," Rev. Mex. Ciencias Pecu., vol. 10, no. 4, pp. 885–902, 2019, doi: 10.22319/rmcp.v10i4.4561.

[16] F. S. Bulitta, Effects of handling on animals welfare during transport and marketing, vol. 2015, no. 117. 2015.

[17] B. L. Nielsen, L. Dybkjr, and M. S. Herskin, "Road transport of farm animals: Effects of journey duration on animal welfare," Animal, vol. 5, no. 3, pp. 415–427, 2011, doi: 10.1017/S1751731110001989.

[18] E. Santurtun and C. J. C. Phillips, "The impact of vehicle motion during transport on animal welfare," Res. Vet. Sci., vol. 100, pp. 303–308, 2015, doi: 10.1016/j.rvsc.2015.03.018.

[19] N. S. Minka and J. O. Ayo, "Physiological responses of food animals to road transportation stress," African J. Biotechnol., vol. 9, no. 40, pp. 6601–6613, 2010.

[20] S. Lee et al., "Body weight changes of laboratory animals during transportation," Asian-Australasian J. Anim. Sci., vol. 25, no. 2, pp. 286–290, 2012, doi: 10.5713/ajas.2011.11227.

[21] E. Panel and W. Ahaw, "Scientific Opinion Concerning the Welfare of Animals during Transport," EFSA J., vol. 9, no. 1, pp. 1–44, 2011, doi: 10.2903/j.efsa.2011.1966.

[22] M. Yerpes, P. Llonch, and X. Manteca, "Effect of environmental conditions during transport on chick weight loss and mortality," Poult. Sci., vol. 100, no. 1, pp. 129–137, 2021, doi: 10.1016/j.psj.2020.10.003.

[23] N. Bachelard, "Animal transport as regulated in Europe: a work in progress as viewed by an NGO," Anim. Front., vol. 12, no. 1, pp. 16–24, 2022, doi: 10.1093/af/vfac010.

[24] Europan Court Of Auditors, Transport of live animals in the EU: challenges and opportunities. 2023.

[25] U. Sara, D. P. Rahardja, H. Sonjaya, and M. Azhar, "Changes in Physiological Condition of Broiler Chickens Sprayed with Water before Transportation," J. Ilmu Ternak dan Vet., 2022, [Online]. Available: https://api.semanticscholar.org/CorpusID:251851124

[26] X. Zhang and J. Zhang, "Analysis and research on library user behavior based on apriori algorithm," Meas. Sensors, vol. 27, no. May, p. 100802, 2023, doi: 10.1016/j.measen.2023.100802.

[27] M. M. Hassan, A. Karim, S. Mollick, S. Azam, E. Ignatious, and A. S. M. F. Al Haque, "An Apriori Algorithm-Based Association Rule Analysis to detect Human Suicidal Behaviour," Procedia Comput. Sci., vol. 219, pp. 1279–1288, 2023, doi: 10.1016/j.procs.2023.01.412.

[28] Hopi Siti Hopipah, Jajam Haerul Jaman, and Ultach Enri, "Web Usage Mining Guna Analisis Pola Akses Pengunjung Website dengan Association Rule," SATIN - Sains dan Teknol. Inf., vol. 7, no. 2, pp. 53–63, 2021, doi: 10.33372/stn.v7i2.735.

[29] M. Dehghani and Z. Yazdanparast, "Discovering the symptom patterns of COVID-19 from recovered and deceased patients using Apriori association rule mining," Informatics Med. Unlocked, vol. 42, no. June, p. 101351, 2023, doi: 10.1016/j.imu.2023.101351.

# Image Caption Generation using Deep Learning For Video Summarization Applications

Mohammed Inayathulla[1], Karthikeyan C[2]

Research Scholar, Department of Computer Science and Engineering[1]
Koneru Lakshmaiah Education Foundation, Green Fields, Vaddeswaram, Guntur, Andhra Pradesh, India[1]
Associate Professor, Department of Computer Science and Engineering[2]
Koneru Lakshmaiah Education Foundation, Green Fields, Vaddeswaram, Guntur, Andhra Pradesh, India[2]

*Abstract*—In the area of video summarization applications, automatic image caption synthesis using deep learning is a promising approach. This methodology utilizes the capabilities of neural networks to autonomously produce detailed textual descriptions for significant frames or instances in a video. Through the examination of visual elements, deep learning models possess the capability to discern and classify objects, scenarios, and actions, hence enabling the generation of coherent and useful captions. This paper presents a novel methodology for generating image captions in the context of video summarizing applications. DenseNet201 architecture is used to extract image features, enabling the effective extraction of comprehensive visual information from keyframes in the videos. In text processing, GloVe embedding, which is pre-trained word vectors that capture semantic associations between words, is employed to efficiently represent textual information. The utilization of these embeddings establishes a fundamental basis for comprehending the contextual variations and semantic significance of words contained within the captions. LSTM models are subsequently utilized to process the GloVe embeddings, facilitating the development of captions that keep coherence, context, and readability. The integration of GloVe embeddings with LSTM models in this study facilitates the effective fusion of visual and textual data, leading to the generation of captions that are both informative and contextually relevant for video summarization. The proposed model significantly enhances the performance by combining the strengths of convolutional neural networks for image analysis and recurrent neural networks for natural language generation. The experimental results demonstrate the effectiveness of the proposed approach in generating informative captions for video summarization, offering a valuable tool for content understanding, retrieval, and recommendation.

*Keywords—Video summarization; deep learning; image caption synthesis; densenet201; GloVe embeddings; LSTM*

## I. INTRODUCTION

Making captions for images is an interesting and useful area of computer vision and natural language processing. The process involves developing algorithms and models that enables machines to provide descriptive and contextually appropriate textual captions for images [24] [25]. This technological advancement serves to connect visual and textual data, so enabling deeper understanding of image content and creating opportunities for diverse applications [1]. The field of image captioning is gaining considerable interest owing to its capacity to boost image accessibility, assist individuals with visual impairments, automate content creation, and enhance image retrieval systems [2]. Especially in video summarization, image caption generation is a potent tool with uses that go beyond individual images. The goal of video summarization is to reduce long videos' main points to more manageable chunks so that viewers may quickly understand the main points without having to watch the full thing. Image caption generation is essential in this situation [3].

The process of video summarization often entails splitting a video into a series of frames, which are effectively separate images. After that, approaches for image captioning are used to these frames, which results in the generation of written descriptions for each frame. There are many ways that image caption creation might be used to the process of video summarization [4]. It may be used in news collection, which provides consumers with the ability to quickly interpret the most important aspects of news broadcasts or events. In educational settings, it can facilitate efficient learning by providing concise summaries of lengthy video lectures. It may be used by content makers to generate appealing video teasers or trailers, and it also has potential applications in surveillance and security, where it might assist analysts in more quickly reviewing video material [5]. The use of image captioning in video summarization not only makes the process of processing material more streamlined, but it also improves searchability and the ability to retrieve certain video portions. This convergence of computer vision and natural language processing puts us one step closer to developing video summarization tools that are more efficient and user-friendly. Due to these challenges, the field of video summarization [6] has seen the emergence of deep learning as a very promising approach. Deep learning, specifically deep neural networks, has an impressive capacity to autonomously acquire complex features and patterns from unprocessed data [27][28]. Furthermore, it can efficiently capture the temporal relationships present in video streams. The process involves the use of sophisticated neural networks and algorithms to examine video streams and identify frames that most effectively depict the information or occurrences inside the film. The frames that have been chosen are often known as "image captures" and function as succinct summaries of the video material [7]. This methodology exhibits a range of realistic implementations, including the retrieval of video material, security and surveillance operations, video search functionalities, and content analysis activities. The primary objective is to discern significant information promptly and effectively inside vast collections of video data. The use of deep learning techniques

enhances the efficacy and adaptability of the summarization process, enabling its application across diverse video kinds and areas [8]. It has an intrinsic capability to modify and generalize over a wide range of video genres, making it an attractive option for tackling the complex challenges presented by video summarization.

The objective of this research is to explore the integration of deep learning methodologies with video summarization, specifically emphasizing the production of picture snapshots. It investigates the capabilities of deep learning models in autonomously selecting and constructing cohesive image-based summaries of movies. It presents a valuable tool with wide-ranging applications in many fields such as surveillance, content analysis, and video search. The subsequent parts provide an overview of the most recent deep learning approaches; proceed throughout their benefits for video summarization, and show case studies and experimental findings that demonstrate how well these methods work to solve the problems associated with creating image captures from video streams. The organization of the paper as follows: Section II describes the literature survey and the proposed model is explained in Section III. The simulation results are discussed in Section IV. Finally, Section V concludes the paper.

## II. RELATED WORK

Hafiz Burhan Ul Haq et al. [9] proposed a deep learning based system for tailored video summarization. The suggested framework facilitates video summarization based on the Object of Interest (OoI), such as individuals, aircraft, mobile devices, bicycles, and automobiles. Sridevi et al. [10] demonstrated with the use of a deep convolutional neural network in each stream, to create a video summary by extracting the temporal and spatial information from a video. The use of a Two-dimensional Convolutional Neural Network (2D CNN) allows for the exploitation of spatial information, while a Three-dimensional Convolutional Neural Network (3D CNN) is employed to exploit temporal information in order to provide highlight scores for video segments. The fusion of segment ratings from each stream is used to identify highlight portions within the video. Moreover, the resulting highlight representation alone indicates the user's relative degree of interest in a movie. To train the deep convolutional neural network (DCNN) in each stream, a paired deep ranking model is used. The objective is to enhance the highlight score of the highlight segment relative to the non-highlight section via the optimization of the model. The segments that have been acquired are then used in the process of summarizing a video.

Obada Issa et al. [11] illustrated novel methodologies for addressing the issue of key frame extraction in the context of video summarization. The methodology used in the study involves the extraction of feature variables from the bit streams of coded films, which is then followed by an optional stepwise regression process aimed at reducing dimensionality. After extracting the features and reducing their dimensionality, novel frame-level temporal subsampling approaches are used, followed by training and testing using deep learning architectures. The frame-level temporal subsampling approaches rely on the use of cosine similarity and the

application of PCA projections on feature vectors. Three distinct learning architectures are constructed by using LSTM networks, 1D-CNN networks, and random forests.

Xu Wang et al. [12] presented a novel deep summarizing network that incorporates auxiliary summarization losses in order to effectively tackle the aforementioned issue. The incorporation of an unsupervised auxiliary summarization loss module using LSTM and a swish activation function is proposed. This module aims to effectively capture long-term dependencies for video summarizing tasks. Furthermore, the proposed module can be seamlessly incorporated into diverse network architectures. The presented model is a novel unsupervised framework for deep reinforcement learning that operates independently of any explicit labels or user interactions. In addition, the suggested model has a low computational burden and may effectively be implemented on mobile devices, hence improving the mobile user experience and alleviating strain on server operations.

Rhevanth et al. [13] presented an effective video summarizing method that extracts essential frames from raw video input and analyzes visual and audio material. Mel-frequency cepstral coefficient (MFCC) extracts information from audio sources, whereas structural similarity index compares frames. Using the preceding two functions removes superfluous video frames. A deep convolution neural network (CNN) model refines the key frames to get a list of potential key frames that summarize the data. Gulraiz Khan et al. [14] suggested a method facilitating users in generating video summaries by using human and object attributes. Cryptographic hashes play a crucial role in the context of blockchain technology. These hashes are derived from condensed video blocks, serving as a means of summarizing the content. Subsequently, these hashes are signed and transferred over the blockchain network. The Cumulus blockchain method is used to safeguard the integrity of the video. The system facilitates distant users in obtaining tamper-proof, condensed video footage of their company locations or other critical properties, which can be accessed on their cellphones. Xiaoning Chen et al. [15] presented a method for video summarization (VS) that leverages the complementary nature of shallow and deep features. The suggested approach involves Multiview feature co-factorization based dictionary selection, which aims to use the shared information from both shallow and deep view features in VS. In order to effectively use the whole visual information of video frames, two view features are employed. Subsequently, the shared information between these two distinct views is extracted using coupled matrix factorization. This retrieved information is then utilized for the purpose of dictionary selection in the context of visual surveillance. Ke Zheng et al. [16] presented a video summarization generation model called DME-VSNet, which utilizes a multi-feature approach to extract various information from the video frames. This study incorporates three key variables: significance score, picture memory strength, and image entropy. In response to the issue of imprecise video shot segmentation, this study presents a video shot segmentation algorithm that utilizes the TransNet network. The system effectively partitions the original video into many shorter shots by identifying shot borders. The suggested model incorporates

three specific variables as inputs for this purpose. The video frame score is acquired inside the Multi-Layer Perceptron (MLP) architecture, and subsequently, the key frame is determined based on this score to provide a concise summary of the movie.

Balamurugan et al. [17] illustrated a model for anomalous event detection that combines a hybrid convolution neural network (CNN) with bi-directional long short term memory (Bi-LSTM). The model is designed to have decreased complexity. The proposed model incorporates a convolutional neural network that utilizes a pre-trained model to extract spatio-temporal features from individual frames within a series. These features are subsequently fed into a multi-layer bi-directional long short-term memory network, which is capable of accurately classifying abnormal events in complex surveillance scenes on the road. The fine-grained technique incorporates a hierarchical temporal attention-based LSTM encoder-decoder model to provide an enhanced video summarizing approach that effectively maintains critical information while optimizing storage capacity.

Sah, Ramesh Kumar et al. [18] proposed a framework that utilizes spatial and temporal aspects, including self-attention mechanisms, to select representative content from video sequences. The framework generates temporal proposals and employs supervised learning techniques using manually provided data from individuals or users. Current supervised approaches do not effectively address the temporal interest and its consistency. In addition, achieving temporal consistency requires the ability to anticipate the temporal suggestions of the video segment. The present study approaches the task as temporal action detection, whereby it aims to concurrently forecast the relevance score and placement of the segments. This is achieved by using an anchor-based system that creates anchors of different lengths to effectively identify intriguing ideas.

## III. PROPOSED METHODOLOGY

A deep learning-based image caption generator is a framework that automatically generates informative text captions for images. By merging computer vision and natural language processing methodologies, this system is capable of comprehending and articulating the semantic content of an image in a manner that is comprehensible to humans.

### A. Image Caption Generator Framework

The image caption generator framework consists of data collection, data presentation, model presentation, and training and validation phases. During the data preprocessing stage, images undergo several operations such as scaling, normalization, and augmentation to ensure uniform dimensions and improved feature representation. Captions/Textual descriptions undergo the process of tokenization, wherein they are segmented into individual units, and subsequently transformed into numerical representations. This conversion is commonly achieved through the utilization of word embedding techniques. The model architecture is implemented using a pre-trained convolutional neural network (CNN), an encoder, which is responsible for processing the image and extracting visual features at a higher level. On the other hand, a decoder,

commonly implemented as a recurrent neural network (RNN), utilizes these extracted features to generate textual descriptions. During the training process, the model aims to reduce the disparity between the captions generated by the model and the actual captions by utilizing a loss function, such as cross-entropy. Fig. 1 shows the Image Caption Generator Framework.



Fig. 1. Image caption generator framework.

*1) Image preprocessing:* Resizing: To maintain consistency, images in the dataset are frequently scaled to a fixed dimension (such as 224x224). The necessity of performing the resizing phase arises from the fact that deep learning models, particularly convolutional neural networks (CNNs), commonly necessitate input images to possess equal dimensions.

Normalization refers to the process of scaling image pixel values to a predetermined range, commonly denoted as [0, 1] or [-1, 1]. Normalization is a crucial step in data preprocessing that aims to establish a uniform range and mean for the input data. This process plays a significant role in enhancing the training process and facilitating the convergence of the model.

Data augmentation strategies are employed in order to enhance the variety of training data and bolster the resilience of the model. This may encompass various image processing processes such as rotation, cropping, flipping, brightness tweaks, and zooming. Data augmentation is a technique that aids in enhancing the generalization capabilities of a model towards images that have not been previously encountered.

*2) Text preprocessing:* Tokenization involves the process of breaking down captions, into individual words or sub words. Tokenization refers to the computational procedure of dividing a given text into distinct and meaningful parts, which might include individual words or sub word tokens. This process is commonly accomplished through the utilization of tools such as the Natural Language Toolkit (NLTK) or spaCy. Every token is representative of either a complete word or a fragment of a word.

Vocabulary creation involves extracting the tokens present in the dataset and organizing them into a comprehensive collection. This lexicon encompasses all distinct lexical units or morphological constituents found inside the captions. The size of the vocabulary is governed by the quantity of distinct tokens. Restricting the size of the vocabulary is crucial to maintain computational efficiency during training and inference, as excessively large vocabularies can lead to increased computational costs.

Padding and sequence length issues arise when dealing with captions, as they frequently vary in length. However, neural networks require input sequences of a set length. Consequently, it is possible to add a specific token (such as <PAD>) to sequences in order to achieve consistent length. The determination of a maximum sequence length is also utilized to appropriately truncate or pad sequences. Captions that are lower than the maximum allowable length are extended by adding padding, and captions that exceed the maximum length are shortened by truncation. Word embeddings are frequently employed to transform words into numerical representations. Pre-existing word embeddings, such as Word2Vec, GloVe, and FastText, can be employed to establish a mapping between words and compact vector representations. Embeddings can capture semantic links between words, hence enhancing the model's comprehension of the text. Special tokens, such as "<START>" to denote the initiation of a series and "<END>" to signify its conclusion, are incorporated into the tokenized captions. The utilization of tokens aids the model in acquiring knowledge regarding the appropriate instances to initiate and conclude the process of generating captions during the decoding phase.

Data preprocessing is a crucial step in preparing picture and text data for effective training of deep learning models. The preprocessed data is subsequently utilized to train the image encoder, which is typically a Convolutional Neural Network (CNN), and the text decoder, which is typically a Recurrent Neural Network (RNN) model, in the image captioning framework. The alignment between the data and the model architecture is crucial for effectively training the model and producing coherent image captions.

*3) Image encoder:* The function of the image encoder is to undertake the processing of the input image and extract significant features from it. Convolutional Neural Networks (CNNs) are frequently employed as image encoders.

*a) Pre-trained CNN:* A pre-trained (CNN) model, such as VGG, ResNet, or Inception, is employed as backbone for the image encoder. These models have already been trained to extract hierarchical and meaningful visual features from images using data from large-scale image classification tasks like as ImageNet.

*b) Feature extraction:* To extract features, the input image is processed through the pre-trained CNN. As the input data traverses the many layers of the CNN, distinct characteristics are identified and represented at varying degrees of abstraction. The properties serve to capture and represent relevant information related to the edges, textures, forms, and constituent components of the depicted image.

*c) Dense layer:* Dense layer or fully connected layer is employed to further transform the feature vector into a feature map that aligns with the desired input size for the text decoder. Embedding of Image Features: The image encoder produces a feature vector as its final output, which serves as a representation of the visual material contained inside the image. The feature vector serves as the initial hidden state for the text decoder. The process of embedding image features into a dense vector is commonly employed to ensure compatibility with the RNN decoder.

*4) Text decoder:* The text decoder reads the visual features and provides textual captions word by word. Text decoders in natural language processing (NLP) often employ Recurrent Neural Networks (RNNs) or Long Short-Term Memory (LSTM) networks. The following is a comprehensive elucidation of the text decoder:

*a) Initial state for the text decoder* is derived from the feature vector obtained from the image encoder. This initializes the decoder by utilizing a reference point to generate captions that are derived from the visual characteristics of the image.

*b) Embedding layer:* The input word tokens or sub word tokens obtained from the preprocessed captions undergo a process of passing via an embedding layer. The process involves the mapping of each word to a dense vector representation, often of a predetermined size. The acquisition of these embeddings takes place during the training process.

*c) Recurrent layers:* The fundamental component of the text decoder consists of the recurrent layers, specifically the Long Short-Term Memory (LSTM) . The layers receive the embedded word representations as input and iteratively update their hidden states. During each iteration, the decoder generates a prediction for the subsequent word in the caption. The hidden state is modified by using information from both the previously created words and the image features.

*d) Output layer:* The output layer of the decoder generates a probability distribution across the vocabulary of words at each time step. The generation of this distribution is accomplished by applying a softmax layer on the hidden state. The subsequent word in the caption is selected based on its highest probability. The integration of an image encoder and text decoder inside a sequential framework constitutes the fundamental component of the image captioning model. The objective of training this model is to reduce the disparity between the captions generated by the model and the ground truth captions obtained from the dataset. The model acquires

the ability to produce coherent and contextually appropriate captions for a diverse array of images.

### B. Proposed Deep Learning Architecture

The CNN-LSTM model that has been presented will be discussed in this section. Fig. 2 illustrates the CNN-LSTM model that has been proposed for use in image captioning. The diagram depicts the various components of the model. The model is made up of several different layers. The CNN layer is used to first extract characteristics from the image. The CNN layer acquires the knowledge necessary to extract features from images, including edges, shapes, and colors, that are crucial to the process of image captioning. After the output of the CNN layer has been formed into a series of vectors, the next layer is the output of the CNN layer. Each vector represents a different part of the image. The sequence of vectors is then passed to the embedding layer. Each vector is inserted into a space with a high dimension by using the embedding layer. This enables the LSTM layer to learn more complex correlations between the vectors. Long-term dependencies in the vector sequence are learned by the LSTM layer. This is crucial for image captioning since a caption should make sense and be in line with the picture's content. The dropout layer receives the LSTM layer's output after that. By removing part of the LSTM layer neurons at random, the dropout layer stops overfitting. The model is compelled to pick up stronger characteristics as a result. Next, the output from the dropout layer is combined with the output from the layer before it. The purpose of doing this is to depict the picture in a more sophisticated way. The model's last layer is a thick layer. The image's caption is produced by the thick layer. The algorithm predicts one word at a time to create the caption. Long Short-Term Memory is an architecture for recurrent neural networks (RNNs)[26] that is meant to manage and analyze sequences of data, such as time series, natural language, and voice. It is abbreviated as LSTM, which is also the name of the corresponding abbreviation. classic RNNs have difficulty collecting long-term dependencies in data, therefore researchers came up with the idea of LSTMs to solve this problem and overcome the limits of classic RNNs. The capacity of LSTMs to successfully acquire and remember information over extended periods is largely responsible for the explosion in popularity of this type of model. The most important innovation of LSTMs is found in their memory cells, which give them the ability to store information and keep it up to date over time. These memory cells are made up of three gates: the input gate, the forget gate, and the output gate. These are the most important gates. The forget gate selects what information is no longer relevant, the input gate regulates what information is stored in the memory cell, and the output gate decides what information is shown to the network's output. These gates are controlled by three gates: the input gate, the forget gate, and the output gate. This architecture not only allows LSTMs to recognize and recall patterns or relationships within sequential data, but it also helps typical RNNs avoid the problem of vanishing gradients, which is a common issue for these types of networks. DenseNet-201 is a convolutional neural network architecture that was devised by Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger in the year 2016. This model is a modification of the DenseNet-121 architecture, aimed at

mitigating the drawbacks commonly observed in conventional deep neural networks like the issue of vanishing gradients and the challenges associated with training very deep networks. Fig. 2 shows the proposed architecture.



Fig. 2. Proposed model.

### IV. RESULTS AND DISCUSSION

This section presents a comprehensive account of the outcomes derived from the simulations carried out utilizing the suggested methodology. The dataset utilized in this research was obtained from Kaggle. The dataset (Flickr 8k) was subjected to processing utilizing the suggested technique. The dataset comprises 8,000 photos, each accompanied by five distinct captions. These captions aim to offer comprehensive descriptions of the prominent things and events depicted in the images. The selection of photographs was derived from six distinct Flickr groups, predominantly devoid of prominent individuals or recognizable landmarks. Fig. 3 shows some sample images and their text captions.

Tokenization is a key approach in the field of natural language processing (NLP) that entails the segmentation of text or language into smaller components known as tokens. The tokens employed in this context generally consist of words, sentences, or even individual characters, depending on the task

and the level of detail desired. Tokenization plays a pivotal role as an essential pre-processing step in numerous natural language processing (NLP) applications, facilitating the comprehension and manipulation of human language by computers. In the domain of English language processing, tokenization conventionally entails the division of words by means of spaces or punctuation marks. However, for languages without distinct word delimiters, the process of tokenization may exhibit greater intricacy. The process of tokenization holds significant importance since it facilitates the study of textual data, including many tasks such as text categorization, sentiment analysis, machine translation, and information retrieval. The utilization of algorithms enables the processing of structured and quantified linguistic data, facilitating the extraction of significant insights from textual information, and facilitating effective communication between humans and machines.

The process of image feature extraction in the field of computer vision entails the conversion of unprocessed picture data into a collection of numerical or symbolic descriptors, commonly referred to as features. These features enable more efficient processing and analysis by machine learning algorithms. The process of feature extraction is of utmost importance as it serves to streamline the intricacies associated with visual data, while preserving the vital information required for a multitude of computer vision applications, including but not limited to object recognition, image classification[22], and image retrieval. The techniques employed for feature extraction exhibit a wide range of complexity, encompassing rudimentary approaches such as color histograms and edge detection, as well as sophisticated methods like convolutional neural networks (CNNs) that possess the ability to autonomously acquire pertinent features from images. The collected characteristics play a fundamental role in picture comprehension and facilitate the development of models that possess the ability to identify objects and detect patterns within images. In brief, the process of image feature extraction serves to connect unprocessed visual input with machine learning algorithms, hence enabling the comprehension and examination of images across several domains, including but not limited to autonomous driving and medical imaging.

The notions of training loss and validation loss hold significant importance in the training and evaluation of machine learning models, specifically in the domain of supervised learning tasks like classification and regression. These metrics serve as important indicators for evaluating the performance of the model both during the training process and in subsequent assessments. The training loss is a metric used to evaluate the performance of a machine learning model on the training dataset. Fig. 4. depicts the training and testing loss of the proposed model. The quantification of the inaccuracy or disparity between the predictions made by a model and the actual target values in the training data is referred to as the evaluation of model performance. Throughout the training procedure, the model iteratively modifies its parameters, such as weights and biases, to minimize the loss function. A decrease in training loss signifies that the model is progressively improving its ability to appropriately match the

training data. Nevertheless, it is imperative to acknowledge that an excessively low training loss does not guarantee that a model would exhibit strong generalization capabilities when presented with unseen data. Indeed, the phenomenon of obtaining a significantly reduced training loss while exhibiting poor generalization performance serves as an indication of overfitting. Overfitting occurs when a model has efficiently memorized the training dataset yet lacks the ability to accurately predict outcomes for novel, unseen data instances.



Fig. 3. Sample images and captions of flickr 8k dataset.

The validation loss, also known as the test loss, is a metric that evaluates the performance of a model on unseen data, which was not used for training. The validation dataset, which comprises unseen data, serves the purpose of evaluating the model's capacity to generalize, and is separate from the training dataset. The calculation of the validation loss is performed similarly to that of the training loss, where the model's predictions are compared against the actual target values. A low validation loss indicates that the model is effectively generalizing its learned patterns to previously unseen data. The metric functions as a significant determinant of a model's efficacy in forecasting and its capacity to generate precise predictions is when applied to the real-world dataset. The results of the proposed model are depicted in Fig. 5.



Fig. 4. Training and test loss of the model.



Fig. 5. Captions generated by proposed model on sample test images.

The BLEU metric [23], known as Bilingual Evaluation Understudy, is commonly employed to assess the caliber of machine-generated text, particularly within the domain of machine translation. The metric quantifies the degree of similarity between the text generated by the machine and a reference text, yielding a numerical score that assesses the alignment of the machine-generated text with the reference text produced by humans. The BLEU metric is commonly employed in the fields of natural language processing and machine translation for evaluating the effectiveness of language generating models. The BLEU score operates by doing a comparison between the n-grams, which are consecutive sequences of n words or characters, present in the machine-generated text and those found in the reference text. The evaluation metric evaluates the precision, which quantifies the ratio of n-grams in the text generated by the machine that are also present in the reference text. Additionally, it considers brevity, which takes into consideration the length of the generated text in relation to the reference text. The BLEU score is quantified as a numerical value ranging from 0 to 1, where a higher score signifies a stronger correspondence between the generated text and the reference text. The BLEU score obtained by proposed method is 0.6052. The BLEU metric of the proposed model is compared with existing models and corresponding results are reported in Table I. The Table I compares models' image caption generation performance using the BLEU score, a typical criterion for machine-generated text quality. These models establish image descriptions using CNN and LSTM combinations. The findings show different performance levels: VGG16, a common CNN architecture, scored 0.56 in BLEU, whereas DensNet+LSTM scored 0.57. The Conventional CNN+LSTM model scored 0.39. The Proposed CNN+LSTM have the greatest BLEU score of 0.60, showing its capacity to create picture captions that resemble human-generated reference captions.

TABLE I. BLEU COMPARISON RESULTS

| S.No | Model | BLEU |
|------|-------|------|
| 1 | VGG16 [19] | 0.56 |
| 2 | DensNet+LSTM [20] | 0.57 |
| 3 | Conventional CNN+LSTM [21] | 0.39 |
| 4 | Proposed CNN+LSTM | 0.60 |

## V. CONCLUSION

In conclusion, this research has showcased the capabilities of deep learning methodologies in generating image captions for video summarization. By utilizing the DenseNet201 architecture for extracting image features and deploying GloVe LSTM models for text processing, the proposed model has effectively developed a framework that effectively connects visual and textual content, providing a comprehensive solution for applications related to video summarization. The captions provided offer significant contextual information and important perspectives for video content, hence enhancing its accessibility and interpretability. The proposed framework obtained a BLEU score of 0.60. The image caption model built using the Flickr8k dataset with proposed architecture has limitations stemming from the dataset's relatively small size,

potentially leading to overfitting and a lack of generalization to diverse images. Access to larger and more diverse datasets, possibly incorporating more specialized datasets for nuanced image understanding, will enhance the model's generalization capabilities.

REFERENCES

[1] Poongodi, M., Mounir Hamdi, and Huihui Wang. "Image and audio caps: automated captioning of background sounds and images using deep learning." *Multimedia Systems* (2022): 1-9.

[2] Dognin, Pierre, Igor Melnyk, Youssef Mroueh, Inkit Padhi, Mattia Rigotti, Jarret Ross, Yair Schiff, Richard A. Young, and Brian Belgodere. "Image captioning as an assistive technology: lessons learned from VizWiz 2020 challenge." *Journal of Artificial Intelligence Research* 73 (2022): 437-459.

[3] Apostolidis, Evlampios, Eleni Adamantidou, Alexandros I. Metsai, Vasileios Mezaris, and Ioannis Patras. "Video summarization using deep neural networks: A survey." *Proceedings of the IEEE* 109, no. 11 (2021): 1838-1863.

[4] Hussain, Tanveer, Khan Muhammad, Weiping Ding, Jaime Lloret, Sung Wook Baik, and Victor Hugo C. de Albuquerque. "A comprehensive survey of multi-view video summarization." *Pattern Recognition* 109 (2021): 107567.

[5] Behrens, Ronny, Natasha Zhang Foutz, Michael Franklin, Jannis Funk, Fernanda Gutierrez-Navratil, Julian Hofmann, and Ulrike Leibfried. "Leveraging analytics to produce compelling and profitable film content." *Journal of Cultural Economics* 45 (2021): 171-211.

[6] Fajtl, Jiri, Hajar Sadeghi Sokeh, Vasileios Argyriou, Dorothy Monekosso, and Paolo Remagnino. "Summarizing videos with attention." In *Computer Vision–ACCV 2018 Workshops: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers 14*, pp. 39-54. Springer International Publishing, 2019.

[7] Dilawari, Aniqa, and Muhammad Usman Ghani Khan. "ASoVS: abstractive summarization of video sequences." *IEEE Access* 7 (2019): 29253-29263.

[8] Tiwari, Vasudha, and Charul Bhatnagar. "A survey of recent work on video summarization: approaches and techniques." *Multimedia Tools and Applications* 80, no. 18 (2021): 27187-27221.

[9] Ul Haq, Hafiz Burhan, Muhammad Asif, Maaz Bin Ahmad, Rehan Ashraf, and Toqeer Mahmood. "An effective video summarization framework based on the object of interest using deep learning." *Mathematical Problems in Engineering* 2022 (2022).

[10] Sridevi, M., and Mayuri Kharde. "Video summarization using highlight detection and pairwise deep ranking model." *Procedia Computer Science* 167 (2020): 1839-1848.

[11] Issa, Obada, and Tamer Shanableh. "Static Video Summarization Using Video Coding Features with Frame-Level Temporal Subsampling and Deep Learning." *Applied Sciences* 13, no. 10 (2023): 6065.

[12] Wang, Xu, Yujie Li, Haoyu Wang, Longzhao Huang, and Shuxue Ding. "A Video Summarization Model Based on Deep Reinforcement Learning with Long-Term Dependency." *Sensors* 22, no. 19 (2022): 7689.

[13] Rhevanth, M., Rashad Ahmed, Vithik Shah, and Biju R. Mohan. "Deep Learning Framework based on audio–visual features for video summarization." In *Advanced Machine Intelligence and Signal Processing*, pp. 229-243. Singapore: Springer Nature Singapore, 2022.

[14] Khan, Gulraiz, Saira Jabeen, Muhammad Zeeshan Khan, Muhammad Usman Ghani Khan, and Razi Iqbal. "Blockchain-enabled deep semantic video-to-video summarization for IoT devices." *Computers & Electrical Engineering* 81 (2020): 106524.

[15] Chen, Xiaoning, Mingyang Ma, Runfeng Yang, and Yong Peng. "Multiview feature co-factorization based dictionary selection for video summarization." *IET Image Processing* (2023).

[16] Zheng, Ke, and Xiangdi Chen. "Research on video summarization method based on convolutional neural network." In *International Conference on Neural Networks, Information, and Communication Engineering (NNICE)*, vol. 12258, pp. 52-56. SPIE, 2022.

[17] Balamurugan, G., and J. Jayabharathy. "An integrated framework for abnormal event detection and video summarization using deep learning." (2022).

[18] Sah, Ramesh Kumar. "Video Summarization using Spatio-Temporal Features by Detecting Representative Content based on Supervised Deep Learning." PhD diss., Pulchowk Campus, 2021.

[19] Sri Neha, V., B. Nikhila, K. Deepika, and T. Subetha. "A Comparative Analysis on Image Caption Generator Using Deep Learning Architecture—ResNet and VGG16." In *Computational Vision and Bio-Inspired Computing: Proceedings of ICCVBIC 2021*, pp. 209-218. Singapore: Springer Singapore, 2022.

[20] Deng, Zhenrong, Zhouqin Jiang, Rushi Lan, Wenming Huang, and Xiaonan Luo. "Image captioning using DenseNet network and adaptive attention." *Signal Processing: Image Communication* 85 (2020): 115836.

[21] Das, Ringki, and Thoudam Doren Singh. "Assamese news image caption generation using attention mechanism." *Multimedia Tools and Applications* 81, no. 7 (2022): 10051-10069.

[22] Inayathulla, M., Karthikeyan, C. (2022). Supervised Deep Learning Approach for Generating Dynamic Summary of the Video. In: Suma, V., Baig, Z., Kolandapalayam Shanmugam, S., Lorenz, P. (eds) Inventive Systems and Control. Lecture Notes in Networks and Systems, vol 436. Springer, Singapore. https://doi.org/10.1007/978-981-19-1012-8_18.

[23] K. Papineni, S. Roukos, T. Ward and W.-J. Zhu, BLEU: A method for automatic evaluation of machine translation, in: Proceedings of the Annual Meeting of the Association for Computational Linguistics, 2002.

[24] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, Lei Zhang, "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering," 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 6077-6086.

[25] J. Wu, T. Chen, H. Wu, Z. Yang, G. Luo and L. Lin, "Fine-Grained Image Captioning With Global-Local Discriminative Objective," in IEEE Transactions on Multimedia, vol. 23,2021, pp. 2413-2427.

[26] Marc Tanti, Albert Gatt, Kenneth P. Camilleri, "What is the Role of Recurrent Neural Networks (RNNs) in an Image Caption Generator?",Proceedings of the 10th International Conference on Natural Language Generation,2017, Pages 51-60.

[27] Amirian, S., Rasheed, K., Taha, T.R., Arabnia, H.R. (2021). "Automatic Generation of Descriptive Titles for Video Clips Using Deep Learning" In: Advances in Artificial Intelligence and Applied Cognitive Computing. Transactions on Computational Science and Computational Intelligence. Springer, Cham. https://doi.org/10.1007/978-3-030-70296-0_2.

[28] Kelvin Xu, Jimmy Lei Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", Proceedings of the 32nd International Conference on Machine Learning, PMLR 2015,37:2048-2057.

# Evolving Adoption of eLearning Tools and Developing Online Courses: A Practical Case Study from Al-Baha University, Saudi Arabia

Hassan Alghamdi[1], Naif Alzahrani[2]

Department of Computer Information Systems, School of CS and IT, Al-Baha University, Al-Baha, Saudi Arabia[1]
Department of Computer Engineering, School of Engineering, Al-Baha University, Al-Baha, Saudi Arabia[2]

*Abstract*—**eLearning or online learning has gained acceptance worldwide, particularly after the Covid-19 pandemic. Although the pandemic has forced the shift towards this learning mode, there is still a continuous need to improve instructors' cognitive and practical competencies to effectively design and deliver online courses. In this paper, a practical case study from Al-Baha University, a Higher Education Institution (HEI) in Saudi Arabia, is presented, showing the development stages of eLearning at the university and how effective utilization of eLearning tools through a structured methodology in a short time, with minimum resources, has helped to improve the teaching and learning experiences for both instructors and students at the university before the pandemic. Various standards and research techniques have been adopted to develop and assess the methodology and its viability of implementation in other higher education institutions. The findings show the methodology's effectiveness and how it helps Al-Baha University smoothly adapt to the online shift at the onset of the pandemic. The methodology is presented to and gained acceptance and recommendation for application in other HEIs in Saudi Arabia from the committee of eLearning and distance education deans in Saudi Universities in March 2023. It also receives the Anthology Middle East award for community engagement in November 2023.**

*Keywords—eLearning; ICT competencies; Higher Education Institutions (HEIs); Learning Management System (LMS)*

## I. INTRODUCTION

A commonly accepted definition of eLearning is challenging to find [1], [2], yet it can be described as the utilization of different modes of technological tools for the purpose of education, whether these tools are web-based, web-capable or web-distributed [3]. Online learning is a term commonly used to describe education that occurs only via the web without face-to-face contact [3]. With various eLearning adopted terms and styles, it is broadly accepted among academic institutions and across all educational fields [4], [5].

More recently, due to the Covid-19 pandemic at the end of 2019, educational institutions had to move from conventional learning to digital and online learning as part of the health crisis management [6]. Although eLearning and its tools and practices were widely accepted around the globe before the pandemic, as previously stated, many Higher Education Institutions (HEIs) and faculty members in many countries were not ready for the shift to online teaching once the

pandemic occurred [7]. The low level of experience among instructors in adopting new teaching formats and the time constraints to make the shift immediately were affecting the readiness of these institutions [8]. This has led many HEIs to speed up the process of adopting eLearning tools without structured long-term planning [9].

The role of teachers in schools and instructors in HEIs during the pandemic was significant to ensure the continuity of teaching and learning [10]. Recent research shows that instructors in HEIs after the pandemic have appreciated blended and online learning, which allows them to combine the strengths of these teaching styles with the face-to-face style [9]. Jelinska et al. [11] pointed out that teachers with prior experience in dealing with online teaching tools were the best at coping with the challenges faced during the transition of the pandemic and were the most engaging. Similarly, Divjak et al. [12] state that instructors who had prior experiences adopting innovative teaching styles, such as flipped classrooms, before the pandemic had more successful experiences with online teaching implementation during the pandemic. However, instructors' adoption of online tools in HEIs before the pandemic was not pervasive [8]. This has significantly affected their experiences with online teaching during the shift to online mode. The focus should then be on continuous investment and enhancement of Information and Communications Technology (ICT) and eLearning competencies of instructors in HEIs to ensure long-term adoption and adaptation of different teaching and learning styles in HEIs. König et al. [13] emphasize the importance of professional development of digital competencies for teachers to adapt to online teaching. Gameil et al. (2023) [14], in their study findings, pointed out that huge investments must be made to train teachers to improve their instructional designing competencies in terms of cognitive and practical perspectives. They highlighted the issue that employing digital learning platforms after the pandemic requires paying more attention to improving the skills of teachers to cope with the radical changes in the teaching landscape.

Similarly, Svetec et al. [15] state that although instructors in HEIs realize the importance of digitalization support, they are not fully aware of the significant role of digital technologies in enhancing teaching and learning experiences. The enhancement of these experiences and the ability to reach a high level of quality in education are common strategic objectives in HEIs that are sometimes misaligned with

available technologies and the actual roles of instructors in these institutions. Jans (2009) [16] also states that teachers and lecturers may have learned and worked with many hardware and software tools; they, however, still require more didactical skills on how to use Learning Management Systems (LMSs) effectively, and longer courses are required to improve their eLearning competencies. According to Jans (2009) [16], eLearning and blended learning tools can only be learned by experimenting with them. Alshammari (2020) [17] also investigated instructors' behavioral intentions to use LMSs in a higher education institution in Saudi Arabia through several factors, such as age, gender, experience, etc., and found a causal relationship between experience and behavioral intention of LMS usage.

This practical case study paper describes a methodology implemented at a higher education institution in Saudi Arabia to elevate the usage of LMS through structured development and delivery of online courses while considering the enhancement of instructors' ICT and eLearning competencies. The methodology is developed following the framework devised by the authors and introduced in the paper [18]. In short, the framework consists of seven iterative stages that assist HEIs in elevating the usage of eLearning technologies by adopting them systematically, ensuring the involvement of related organizational components from alignment and Enterprise Architecture (EA) perspectives.

The following sections in this paper present background information about eLearning in Saudi Arabia in general and at Al-Baha University in particular. Then, the problems that motivate developing the structured approach to elevate the usage of eLearning technologies and develop and deliver online courses are revealed following the framework of Alzahrani et al. (2023) [18]. After that, the approach and its application and findings are described and discussed. Final thoughts and future work are stated by the end of this paper.

*A. eLearning in Saudi Arabia*

One of the fundamental characteristics of the HE system in Saudi Arabia is that it is centralized in terms of control and educational support [19]. This means that the Ministry of Education (MoE) and several specialized quality assurance and assessment centers are responsible for planning, coordinating, and supervising the kingdom's higher education sector. However, in recent years, the HE system has undergone tremendous reforms, unifying all efforts toward achieving Saudi Vision 2030 [20].

During the last two decades, the government of Saudi Arabia has invested enormously in ICT and digital transformation. According to the Digital Government Authority [21], the cumulative investment size in digital transformation in Saudi Arabia from 2022 to 2025 is expected to reach 25 billion USD. These massive investments in digital technologies contribute towards achieving the government's strategic plan "Vision 2030", where Saudi Arabia aims to be a global technology hub. According to published statistics from [21], Saudi Arabia was ranked second worldwide in digital competitiveness by the European Centre, ranked eleventh among the G20 for digital governance by the Japanese WASEDA ranking, and ranked 1st in the ESCWA indicator for

the availability of digital services, and 19 agencies were awarded Enterprise Architecture Certification from the Digital Government Authority. Proudly, Al-Baha University was among these agencies. There is no doubt that Saudi Arabia has achieved these numbers regarding digital transformation and digital capabilities since it is the largest ICT market in the Middle East and North Africa (MENA) region [22].

The education sector, in particular, receives the largest segment of the government expenditure [22]. Statistics from the same authority show a 100% digital transformation of teaching and learning activities during the pandemic in schools and HEIs through available platforms of "Madrasti" and LMS. Over 350 million virtual classrooms were created during the pandemic, and similarly, large numbers for other activities, such as homework, discussion boards, online exams, etc., were achieved. All have been achieved through a cooperative approach from involved stakeholders and under direct supervision from the Ministry of Education (MoE) and according to the regulations and standards of the National eLearning Centre (NELC).

Abdulrahim et al. [22] state that research on online learning in Saudi Arabia emphasizes the importance of having a strategic direction for its implementation and ensuring that instructors are trained to design and deliver online teaching effectively. Moreover, course design and technology gaps are among the main limitations that face adopting digital learning in Saudi Arabia [23]. Similarly, Alharbi [24] indicates that although there is a broad acceptance of the usefulness of digital learning in Saudi Arabia, there is a lack of knowledge and skills among students and instructors.

*B. eLearning at Al-Baha University*

Al-Baha University is now 17 years old and was established in late 2006. It has 16 faculties spanning six different geographical areas covering the vast region of Al-Baha province (11,000 square km) southwest of Saudi Arabia. It has a current enrolment of 17748 students, 1536 instructors and 776 employees. The university offers 28 undergraduate programs, 17 master programs, and nine programs for higher diplomas. The number of students is increasing annually from only 13,793 in 2008 to nearly 18000 in 2023.

Colbran et al. (2013) [25] studied the adoption of educational technologies among Saudi universities. They found the lack of IT infrastructure, training, and support, as well as websites and software problems, to be the main inhibitors of the successful adoption of educational technologies in the country. However, Saudi Universities have been transforming dramatically in adopting digital technologies along with the enormous investments and support from the government, especially during the last five to eight years. For example, since 2016, Al-Baha University has undergone tremendous changes in adopting and embracing digital technologies and providing services to its key stakeholders, namely students and instructors. By late 2018, the transformation had extended to reform the utilization of eLearning tools at the university after conducting a gap analysis and being aware of the challenges as described next.

The following sections are arranged in the following manner: First, the issues that drive the development of structured approaches to improve the use of eLearning technologies and the creation and delivery of online courses are discussed. Then, the approach is explained, along with its application and findings. Finally, concluding thoughts and future work are presented.

## II. AWARENESS OF THE PROBLEM

Following the proposed framework of Alzahrani et al. (2023) [18], it states that being aware of the challenges that face underutilization of eLearning technologies and linking them with the strategic vision and objectives of HEIs are crucial steps towards effective and continuous usage of these technologies by involved stakeholders. It also suggests having an ecosystem to drive eLearning projects to realize strategic goals and assist in attaining them.

Al-Baha University -similar to other public universities in Saudi Arabia- was given a quota of thirty thousand licenses from the Ministry of Education (MoE) to access an LMS in 2012. This number of licenses is renewed every five years through a national agreement initiative to facilitate and spread the adoption of eLearning tools among students and instructors in these universities.

Although the number of students and instructors at Al-Baha University has never reached this number, the MOE was generous in providing more licenses for any expected future growth in the number of users. Surprisingly, the number of active users using the LMS was less than 0.5%, including students and instructors. This number was revealed in the middle of a running semester in December 2018. The number of registered students and instructors was close to 20000, interacting through 6500 study sections. This number was considered a wake-up call for the university and its management to start a new era of adopting and utilizing eLearning tools effectively.

A gap analysis was conducted to examine early adopted approaches that the university has carried out to embrace the adoption of eLearning tools. The outcomes revealed that the university has:

*1) Fifty-five* developed Electronic Courses (eCourses) that were invested in and transformed from conventional courses to eCourses from 2013 to 2015, following instructors' preferences and available content. Some of these developed eCourses are used only once and never used afterward. Others are used two or three times and similarly put on the shelf. This happens due to the absence of development standards and the intensive customizations applied to these eCourses. Like many other universities in the region during the early 2000s, converting traditional courses to eCourses was the mainstream regardless of its drawbacks [26].

*2) Very* limited attempts to adopt eLearning tools with currently available LMS.

*3) High* resistance to change among students and instructors with negative attitudes towards eLearning.

*4) Many* training sessions and workshops were conducted to explain the functionalities and tools of the LMS, but the follow-up was extremely weak, so what has been learned usually goes missing with less practice in place.

Taking these outcomes into consideration, an eLearning ecosystem named 'Rafid' is created, which consists of a set of components that are linked to the strategic goals of Al-Baha University, aiming to have a holistic transformation and better utilization of eLearning tools at the university.

## III. RAFID ECOSYSTEM

"Rafid" is an Arabic word translated as "assistant/ helper/ supporter." It appeared everywhere at Al-Baha University in December 2018 in a ceremony attended and supported by the university's top management. It was a flag of change in students' and instructors' experiences about online teaching and learning. The word Rafid was used to indicate an eLearning ecosystem and change management approach that involves LMS, virtual labs, analytics tools, methodology for developing online courses, etc., that were all made available to assist, help, and support teaching and learning activities at the university (see Fig. 1). The message was clear to the main stakeholders of the university that Rafid ecosystem is found to assist you in your journey of teaching and learning at the university.

One of Rafid's ecosystem components is the Course Development Lifecycle (CDLS), which, in short, is a methodology that takes instructors on a journey of developing and delivering an online course through a structured approach. It was specifically designed to change the mindset of students and instructors towards embracing eLearning tools through teaching and learning processes.

Four stages are defined in the CDLC. These stages are presented in Fig. 2.



Fig. 1. Rafid ecosystem at Al-Baha university.

Fig. 2.    Course Development Lifecycle (CDLC).

In the first stage, instructors are selected to participate in the process of CDLC based on a nomination process carried out by their faculties according to their ability to allocate 20 hours per week for this project. It was also recommended to nominate instructors who have preliminary ICT competencies such as accessing and browsing the Internet and using some of the main Microsoft Office programs such as Word and PowerPoint to some extent. Once all candidates are nominated, teamwork from the Deanship of eLearning and IT (DeIT), which specializes in course designing and multimedia development, is assigned. Systems and labs are also prepared, the project plan and timeline for course development and delivery are then set, and this stage ends up with a Service Level Agreement (SLA) contract signed by the DeIT from one side and the candidate (i.e., the instructors) from the other side to ensure the commitment of both sides.

In the second stage, participants develop their online courses following best practices and international standards such as Quality Matters (QM) standards for course design. Participants and assigned team members certified in online course design and delivery carry out this stage step-by-step. Assigned team members provide continuous feedback to participants to enhance every aspect of their online courses. This stage focuses on designing and enriching course contents with various resources, creating learning objects, assessment pools, and activities that ensure the engagement and collaboration of students in a standardized and effective learning environment. A master course template (i.e., an online course template structured according to best practices and standards of designing online courses made available for participants on LMS to have hands-on experience while developing their courses step-by-step). In other words, it facilitates developing standardized online courses through what is being learned from the methodology and the assigned certified team member. This course is then made available to students once the development and assessment of fulfilling the requirements are achieved. Several master courses are created and customized according to the nature of the content each faculty delivers at the university.

Once the course is developed, it is revealed to associated students, and the third stage, course delivery, begins. As the name of this stage indicates its objective, the participants deliver their developed online courses to students according to best practices that ensure the instructors' effective, coordinated, and assistant role in delivering online courses.

In the fourth and last stage, an assessment of previous stages is conducted by participants and students. Questionnaires and one-to-one interviews are carried out to find out the strengths and weaknesses of the project and the proposed methodology, more specifically.

Part of the methodology is an online course specifically developed and made available to participants on the LMS to go through the whole experience the students will receive once their online courses are available. This course also helps manage the entire project and lay down the foundations of eLearning and its tools by providing the required materials. Tasks mandatory for participants to complete are also done through assignment tools in the LMS. The project also activates tools such as wikis, discussion boards, groups, virtual classes, etc., where only questions, inquiries, and support needs are accepted and handled through these tools. This is to embrace usability and to show real-life utilization of such tools.

Online (must-attend) workshops provided to participants during the project include forming accurate and measurable learning objectives, designing roadmaps, enhancing and enriching course content with various multimedia content, creating question bank(s), importing and exporting questions, managing tests and assignments and finally delivering effective online courses.

Advanced (optional) online workshops are also offered to participants to expose them to different available techniques and tools that can help them streamline their course design and delivery. Examples of provided workshops include how to use video editing tools and how to build effective presentations through Adobe Suite, Camtasia, Xmind, MindMeister, Microsoft PowerPoint, Prezi, Canva, Kahoot, SurveyMonkey, Mentimeter, Google Docs, and Microsoft OneDrive.

## IV.    APPLYING CDLC METHODOLOGY

In one academic year, 99 instructors have participated and effectively experienced how to develop and deliver online courses following international standards in the field. This was done during two academic semesters in 2019, starting with 24 instructors and then moving to another 75. Participants were from 15 faculties and across different disciplines.

Upon completing the development and delivery of their online courses, they received recognition and certificates from the university president in a rewarding ceremony. They are then recognized as "Rafid ambassadors" in the change management process. They were passionate about passing the message to other instructors and voluntarily engaged in providing support to others in their free time without hesitation, especially during the Covid-19 pandemic.

Similarly, several students were trained to use LMS effectively as student users. We named them "Rafid Friends" and they also started spreading the message of how easy it is to use LMS and its tools to support the learning process at the university. They then provided help and support to other students whenever needed. These students voluntarily shoot videos of themselves using the system, designing brochures and infographics to spread awareness of the LMS and its tools. All these approaches and techniques were undertaken to change the mindset towards eLearning and its tools among students and instructors at the university.

## V.    OUTCOMES AND FINDINGS

Ninety-nine instructors who participated in this project were asked and thankfully completed an online open-ended

questionnaire and attended a 30-minute online interview each. These research tools were designed to receive feedback from the participants about the following:

*1) An* overview of the project in general and the level of satisfaction with the components and the defined stages of the developed methodology used in the CDLS project. This also measures the tendency to reengage in the same experiment again and to recommend it to other instructors, how the whole experiment was useful to them, and how it impacted the teaching and learning practices for them and their students. The outcomes are shown in Fig. 3 and measured on a 5-point Likert scale, with five being very satisfied and one very unsatisfied.

*2) Participants* have learned 46 skills during the project, which enhanced their ICT and eLearning competencies by adopting eLearning tools during the design and delivery processes undertaken throughout the project. The enhancement of ICT and eLearning competencies is not limited to these skills. However, they are drawn from best practices and international standards such as Quality Matters. From the participants ' perspective, the importance of these skills is measured as a percentage, with 100% being very important. Fig. 4 shows the outcomes of this assessment (see Table AI in Appendix A for details of the skills).

*3) The* ability of participants to implement the learned skills mentioned in the previous point. Participants were asked to provide feedback about the easiness of implementing the defined set of skills. They all agree that they can implement each of the mentioned skills with a percentage close to 100%., except (providing online classes through the LMS to students, recording online classes that are being delivered through LMS

and making them available to students, allowing students to communicate through online classes that are created for them, 75%, 70, 60%) accordingly.

*4) The* ability of participants to teach/support others to carry out similar tasks. Participants confirm that they can teach/help others carry out similar learned skills with a percentage exceeding 95% for each skill.

The findings of the interviews conducted with the participants are illustrated in Table I, which outlines the advantages, disadvantages, and recommendations for improving the project of designing and delivering online courses through a developed approach.

A designed open-ended questionnaire is inserted in the developed ninety-nine courses after obtaining permission from the instructors who deliver them, and students who are registered in these courses are asked to fill out this questionnaire. The questionnaire consists of 40 questions asking students about their new eLearning experience. They were given three options against each question as to whether they agree to a great extent, to some extent, or a little extent. A total of 1800 responses were collected from students, and the outcomes are shown in Fig. 5 (see Table BI in Appendix B for details of the skills).

Students are also allowed to express their attitudes about the eLearning experience they have engaged in throughout the project through three open questions at the end of the questionnaire. Questions are focused on gathering feedback about the positive and negative aspects of the experiment as well as any recommendations for future improvements. Responses in each category exceed 100, but the most repeated answers are illustrated in Table II.



Fig. 3. Feedback about the project in general from ninety-nine participants.

Fig. 4. The importance of each learned skill by participants from their perspective.

TABLE I.    ADVANTAGES, DISADVANTAGES, AND RECOMMENDATIONS FOR IMPROVEMENTS FROM PARTICIPANTS

| Advantages | Disadvantages |
|---|---|
| - Clarity of the training course and the flow of steps in each defined phase.<br>- Quick responses and problem-solving from the assigned team members.<br>- Accuracy of follow-up and assessment forms.<br>- Quality of provided training and trainers.<br>- The master course was generally excellent, saving participants time and effort in developing required courses and motivating some participants to continue with the project.<br>- Realizing the benefit of having a road map for each chapter/unit to guide students in the learning flow of the chapter/unit. | - Difficulty engaging in the project and its assignments at the beginning.<br>- Difficulty attending courses and workshops with many assignments while having other administrative roles at the university.<br>- Lack of sufficient labs to conduct online exams on some campuses.<br>- Sometimes, some campuses have Internet and Wi-Fi coverage weaknesses.<br>- There was not sufficient training for students.<br>- Lack of visual content in the training course.<br>- Insufficient time for required assignments.<br>- Sometimes, it was challenging to deal with and modify the content inside the master template in the LMS, especially tables. |
| **Recommendations for improvements** | |
| - Considering the peculiarity of engineering materials.<br>- The project tasks should begin a long time before the semester commences.<br>- A separate e-exam system from LMS should be afforded.<br>- Strengthening the training course with more videos so that it is more of a self-learning course.<br>- Providing participants with a timetable that defines all steps that are required.<br>- There should also be training for students on how effectively they can engage in online learning.<br>- Allow more time for participants to submit the required tasks.<br>- Providing access to virtual experiments and virtual labs in science courses. | |

Fig. 5.   Students' attitudes about the eLearning experience through designed and delivered online courses.

TABLE II.      POSITIVE AND NEGATIVE ASPECTS AND RECOMMENDATIONS FOR IMPROVEMENTS FROM STUDENTS' PERSPECTIVE

| Positive Aspects | Negative Aspects |
|---|---|
| - Easy access to course materials at anytime and anywhere.<br>- Learning has become more flexible.<br>- Easy access to assessments and grades.<br>- Ease of communication and interaction with the course lecturer.<br>- Faster in terms of understanding and accessing information.<br>- Availability of multiple learning materials that support the main subject of study.<br>- Learning through the LMS is fun and motivating for students, which increases students' interest in receiving more information. | - Not all instructors use it.<br>- The number of activities has increased exponentially.<br>- Technical problems sometimes occur, such as (poor Internet connectivity).<br>- Missing assignments and exam submission dates through the LMS due to less familiarity.<br>- Lack of clarity regarding some assessments, such as (time for submitting assignments, grades, etc.).<br>- Most online exams are based on objective questions, and there is no opportunity to express answers that the instructors can understand.<br>- There are no alert notifications in the mobile application that are linked to the LMS.<br>- Not all functionalities are supported or available in the LMS smartphone application. |
| Recommendations for future improvements | |

- Enhancing the mobile application version of the LMS.
- Forcing all instructors to use it in all courses at the university.
- Provide students with some training, specifically regarding online exams.
- Prepare locations with strong Internet connectivity for students to use the system effectively at the university campus.
- Providing technical support 24/7.
- Provide more training to some instructors to use the LMS effectively with students.

Different lessons were learned from the feedback of both students and instructors, and new opportunities were opened to enhance the project in general. Still, this project's most valuable outcome was spreading awareness of eLearning and its tools among students and instructors at Al-Baha University. From only 136 active users and 88 active study sections on the LMS in December 2018, the numbers have increased dramatically to reach 3573 users and 2888 study sections, respectively, after one academic year. This growing number of users represents 35% of the total number of users of the system and 44% of study sections that should be active on the LMS system at that time.

We have utilized the feedback received from this project and have transitioned to a higher-level project where we began designing and providing online courses based on an entire study program level. In this project, an agreement is signed with the dean of selected faculties to have a partnership project targeting the development and delivery of an entire program through their instructors and with the help and support of the specialized team from DeIT. This project confirms the importance of having an iterative approach in embracing eLearning technologies in HEIs, which is suggested in the framework of Alzahrani et al. (2023) [18]. Through this extended and enhanced project, 84 online courses were developed with the faculties of CS and IT (department of Computer Science) and the faculty of Business Administration (department of Marketing), 40 and 44 online courses, respectively. This also helped the growth of the number of users of the LMS, reaching 10163 active users and 3620 active study sections, representing 47% of total expected users and 55% of total active study sections at the university at that time. These numbers were reached before the pandemic of Covid-19 occurred, which helped the university enormously to cope with the healthcare restrictions and the compulsory shift to fully online learning. Instructors who have participated in this project and have been trained very well to use eLearning tools have become trainers to support others in effectively using eLearning tools during the pandemic. We have recorded 210 activities performed by 75 trained instructors during the early months of the pandemic. These activities range from providing videos on how to use eLearning tools to reaching a neighbor during the lockdown period and assisting them in setting up an online exam on the LMS.

Besides enhancing the ICT and eLearning competencies of the instructors at the university, which can be utilized to develop more courses, the developed courses are reusable. 6 of the 24 developed courses in the first semester are reused by the instructors who developed them. Seven other developed courses are voluntarily made available and reused by other instructors. In addition, more than 60 virtual classes were delivered by instructors during the first semester of the project, sending a message to other instructors throughout the university that delivering lectures can be done differently outside the university boundaries.

Before this project, male instructors used to travel to female students' campuses to deliver lectures through video conferencing dedicated rooms. These rooms are entirely locked now and have never been used after the pandemic due to the effectiveness of virtual classes through the LMS. During the first semester of this project, it was also the first time official mid and final exams were carried out electronically through the LMS. This has never happened throughout the history of the university.

There was also a significant increase in the use of the university's official emails among students and instructors, along with the implementation of the project. This is because students receive emails about any new announcements from the LMS, and instructors were also advised to use official emails in all correspondence with students. The number of active email users at the university once the project started was only 2032. This number doubled in one semester only, reaching 3880 active email users.

It is important to note that some recent research argues that attending traditional in-person classes leads to better academic performance than taking online courses [27], [28]. Nevertheless, eLearning tools have transformed the education industry by significantly benefiting students and instructors. Adopting these tools has allowed learners to access information anytime, anywhere, and at their own pace. Meanwhile, instructors have been able to deliver interactive and engaging content to their students from the comfort of their homes. Although eLearning has many benefits, we firmly believe that in-person classes supported by advanced eLearning technologies offer the best possible education now and in the future. This is also supported by recent findings of Chen et al. (2023) [29].

## VI. FINAL REMARKS

The perception of eLearning among Al-Baha University's students and instructors has been positively transformed through the implementation of this comprehensive project, as evidenced by increasing numbers. This project aims to have an effective holistic alignment of eLearning technologies with the strategic scene of Al-Baha University that seeks to enhance the quality of teaching and learning through ICT adoption. This is supported by the framework the authors suggest in Alzahrani et al. (2023) [18]. Although the quality of online courses delivered does not measure the success of this project, it has changed perceptions towards eLearning among critical stakeholders at the university. With the great support of the university's top management, the project was delivered as planned. However, the quality of delivered online courses is iteratively measured by the end of each project implementation cycle as a long-term plan for course evaluation and enhancement.

This project aims to promote the adoption of eLearning in higher education, with a particular focus on students and instructors, while considering other groups such as employees, senior individuals, and different age groups for future projects. The project discussed and the elaborated case study can benefit academic institutions seeking to engage instructors in onrline course development and delivery processes with minimal costs. All participants in the project at Al-Baha University have shown enthusiasm and a sense of responsibility for enhancing their students' learning and teaching experiences. Some members of Saudi universities' eLearning and distance education deans committee suggest implementing a monetary

reward system for project participants to achieve better results, especially after the pandemic.

## VII. Future Work

While implementing our CDLC methodology during the Covid-19 pandemic, we observed the unique aspects of course design and delivery in the faculty of medicine. This is expected as the field of medicine requires specific attention due to the application of various teaching methods and the practical nature of the field. As a result, we studied ten different teaching methods and devised an approach for each method to shift it from face-to-face delivery to either blended hybrid or fully online teaching. We will present and discuss our findings in future work.

## References

[1] V. Arkorful and N. Abaidoo, "The role of e-learning, advantages and disadvantages of its adoption in higher education," International journal of instructional technology and distance learning, vol. 12, no. 1, pp. 29–42, 2015.

[2] A. Ziegler, T. Peisl, and P. Harte, "Linking innovation and eLearning–The case for an embedded design," in European Conference on Software Process Improvement, Springer, 2021, pp. 47–63.

[3] M. Nichols, "A theory for eLearning," J Educ Techno Soc, vol. 6, no. 2, pp. 1–10, 2003.

[4] H. Rodrigues, F. Almeida, V. Figueiredo, and S. L. Lopes, "Tracking e-learning through published papers: A systematic review," Comput Educ, vol. 136, pp. 87–98, Jul. 2019, doi: 10.1016/j.compedu.2019.03.007.

[5] B. Al Kurdi, M. Alshurideh, S. Salloum, Z. Obeidat, and R. Al-dweeri, "An empirical investigation into examination of factors influencing university students' behavior towards elearning acceptance using SEM approach," 2020.

[6] M. Habes, M. Elareshi, E. Youssef, S. Ali, and M. Qudah, "Social impact of videos at new media platforms on the eLearning acceptance during the Covid-19," Inf. Sci. Lett, vol. 11, no. 3, pp. 913–923, 2022.

[7] Y. V. Lakshmi, "eLearning Readiness of Higher Education Faculty Members," Indian Journal of Educational Technology, vol. 3, no. 2, 2021.

[8] P. Santiago, L. Troy, and T. Weko, "The state of higher education: One year into the COVID-19 pandemic," 2021.

[9] A. M. Müller, C. Goh, L. Z. Lim, and X. Gao, "Covid-19 emergency elearning and beyond: Experiences and perspectives of university educators," Educ Sci (Basel), vol. 11, no. 1, p. 19, 2021.

[10] C. Europeo, "Council conclusions on 'European teachers and trainers for the future.'" Retrieved from eur-lex. europa. eu/legal-content/EN/TXT/HTML, 2020.

[11] M. Jelinska and M. B. Paradowski, "Teachers' Engagement in and Coping with Emergency Remote Instruction during COVID-19-Induced School Closures: A Multinational Contextual Perspective.," Online Learning, vol. 25, no. 1, pp. 303–328, 2021.

[12] B. Divjak, B. Rienties, F. Iniesto, P. Vondra, and M. Žižak, "Flipped classrooms in higher education during the COVID-19 pandemic: Findings and future research recommendations," International Journal of Educational Technology in Higher Education, vol. 19, no. 1, pp. 1–24, 2022.

[13] J. König, D. J. Jäger-Biela, and N. Glutsch, "Adapting to online teaching during COVID-19 school closure: teacher education and teacher competence effects among early career teachers in Germany," European journal of teacher education, vol. 43, no. 4, pp. 608–622, 2020.

[14] A. A. Gameil and A. M. Al-Abdullatif, "Using Digital Learning Platforms to Enhance the Instructional Design Competencies and Learning Engagement of Preservice Teachers," Educ Sci (Basel), vol. 13, no. 4, p. 334, 2023.

[15] B. Svetec, L. Oksanen, B. Divjak, and D. Horvat, "Digital teaching in higher education during the pandemic: Experiences in four countries," in Central European Conference on Information and Intelligent Systems, Faculty of Organization and Informatics Varazdin, 2022, pp. 215–222.

[16] S. Jans, "E-learning competencies for teachers in secondary and higher education," International Journal of Emerging Technologies in Learning (iJET), vol. 4, no. 2, pp. 58–60, 2009.

[17] M. H. Alshammari, "Investigating the faculty behavioral intentions to adopt Learning Management Systems (LMSs) in a higher education institution in Saudi Arabia," Virginia Tech, Virginia , 2020.

[18] N. Alzahrani and H. Alghamdi, "Towards a Framework for Elevating the Usage of eLearning Technologies in Higher Education Institutions," International Journal of Advanced Computer Science and Applications, vol. 14, no. 12, 2023, doi: 10.14569/IJACSA.2023.0141230.

[19] L. Smith and A. Abouammoh, "Higher Education in Saudi Arabia," Netherlands: Springer, 2013.

[20] N. Al-Otaibi, "Vision 2030: Religious education reform in the Kingdom of Saudi Arabia," King Faisal Center for Research and Islamic Studies, vol. 1, no. 1, pp. 7–8, 2020.

[21] DGA, "Indicators Of The Digitial Government Authority," https://dga.gov.sa/en/node. Accessed: Sep. 08, 2023. [Online]. Available: https://dga.gov.sa/en/node.

[22] H. Abdulrahim and F. Mabrouk, "COVID-19 and the digital transformation of Saudi higher education.," Asian Journal of Distance Education, vol. 15, no. 1, pp. 291–306, 2020.

[23] A. Aljaber, "E-learning policy in Saudi Arabia: Challenges and successes," Res Comp Int Educ, vol. 13, no. 1, pp. 176–194, 2018.

[24] A. Alharbi, "E-learning in the KSA: A taxonomy of learning methods in Saudi Arabia," 2013.

[25] S. Colbran and N. Al-Ghreimil, "The role of information technology in supporting quality teaching and learning," in Higher education in Saudi Arabia: Achievements, challenges and opportunities, Springer, 2013, pp. 73–82.

[26] J. A. Itmazi and A. S. Alaamer, "Role, objectives, requirements and activities of eLearning units in traditional universities," International Journal of Vocational and Technical Education, vol. 3, no. 3, pp. 29–35, 2011.

[27] H. Marriott, "Can Online Classes Match the Quality of In Person Computer Science Classes?," in 2021 ASEE Pacific Southwest Conference-" Pushing Past Pandemic Pedagogy: Learning from Disruption", 2021.

[28] C. Cardonha, D. Bergman, and R. Day, "Maximizing student opportunities for in-person classes under pandemic capacity reductions," Decis Support Syst, vol. 154, p. 113697, 2022.

[29] Y. Chen and K. Yamamoto, "A study on the Impact of In-Person and Online Formats of University Physical Education Classes on the Acquisition of Life Skills among University Students: Using Japanese University Students as an Example," Advances in Physical Education, vol. 13, no. 3, pp. 151–163, 2023.

APPENDIX A

TABLE AI.     LIST OF LEARNED SKILLS BY PARTICIPANTS

| 1 | Creating a welcome page to students in the LMS |
|---|---|
| 2 | Creating a course description page |
| 3 | Creating a course tour guide inside the course page in the LMS |
| 4 | Creating a forum for knowing each other's activity |
| 5 | Filling out the details of the course description in the LMS |
| 6 | Filling out lecturer's information who is delivering the course |
| 7 | Filling out academic calendar and course tasks' calendar |
| 8 | Filling out course references and materials |
| 9 | Filling out marks policy for the course |
| 10 | Filling out communication policy page |
| 11 | Filling out cyber security policies and requirements requested by cyber security department |
| 12 | Splitting the course in the LMS into chapters/units |
| 13 | Each chapter/unit must have a number of measurable learning objectives |
| 14 | Each chapter/unit has a learning road map |
| 15 | For each course content added, it is linked with a chapter/unit objective |
| 16 | There is a variety of content in each chapter/unit |
| 17 | There is a learning activity in each learning chapter/unit |
| 18 | There is a clear objective for each learning activity |
| 19 | There is a variety of activities throughout the course chapters/units |
| 20 | There is a clear description for each added activity |
| 21 | For each activity, there is a clear criteria for assessment |
| 22 | There is an assignment in each chapter/unit |
| 23 | There is at least one online exam in the course |
| 24 | For each online exam or assignment there must be specified objective(s) |
| 25 | Exams and assignment are linked to course objectives |
| 26 | There is a variety of exams and assignments throughout the course |
| 27 | There is a clear description for each assignment or exams included |
| 28 | For each exam or assignment there must be a clear criteria for assessment |
| 29 | There is a weekly announcement shared with students through LMS |
| 30 | Announcements are used to communicate the details and updates about the course |
| 31 | Different channels are used to communicate with students such as LMS and emails |
| 32 | Discussions in the discussion forums are closely monitored |
| 33 | Students' questions in the discussion forums are answered in acceptable range of time (24hrs max) |
| 34 | Questions and assumptions and examples are used in discussion forums |
| 35 | Positive and negative responses from students in the discussion forums are monitored |
| 36 | Students are encouraged to reply to their classmates' questions in the discussion forums |
| 37 | Responses and discussions in the discussion forums are summarized |
| 38 | Providing online classes through the LMS to students |
| 39 | Allowing students to communicate through online classes that are created for them |
| 40 | Recording online classes that are being delivered through LMS and making them available to students |
| 41 | Ensuring that an effective eLearning environment is maintained throughout the semester |
| 42 | Monitoring students' attendance and participation |
| 43 | Sending regular notifications to students through LMS and emails |
| 44 | Using plagiarism checker tool provided in the LMS |
| 45 | Collecting students feedback about the developed and delivered course |
| 46 | Studying students' feedback about designed and delivered course |

APPENDIX B

TABLE BI.    STUDENTS' ATTITUDES ABOUT ELEARNING EXPERIENCE THROUGH DESIGNED AND DELIVERED COURSES

| | |
|---|---|
| 1 | I always follow course announcements and notifications through LMS |
| 2 | I know about the course description through LMS |
| 3 | I know about assessment methods and policies through LMS |
| 4 | I get to know the course lecturer through LMS |
| 5 | I know required activities and assignments in advance and submit them through LMS |
| 6 | I know my grades from the LMS |
| 7 | I found extra materials provided in the LMS very useful |
| 8 | I can discuss with my lecturer about any topic in the course through LMS |
| 9 | I can share documents with my classmates through LMS |
| 10 | I found email service through the LMS very useful |
| 11 | I can express my opinion about the course and the lecturer through questionnaires made available through LMS |
| 12 | I attend some online lectures in a synchronous mode through the LMS |
| 13 | I do online exams through the LMS |
| 14 | I struggle sometimes to access LMS due to poor Internet connectivity |
| 15 | There are no designated areas or labs to use LMS inside campuses |
| 16 | My ICT competencies are not good enough to use the LMS and its tools |
| 17 | Trainings to use LMS and its tools for students are rarely provided |
| 18 | There is shortage in the awareness about the usefulness of LMS and its tools |
| 19 | There is shortage in providing technical support to students when needed |
| 20 | It is difficult to use and deal with the LMS and its tools |
| 21 | I believe that the LMS is not sufficient enough in terms of its capabilities |
| 22 | Instructors are not prepared or good enough to use the LMS |
| 23 | Delivered materials through the LMS are with low quality |
| 24 | I believe that learning through LMS makes the learning process much easier |
| 25 | I do not want to attend training courses that show me how to use the LMS effectively |
| 26 | I enjoyed learning my courses through the LMS |
| 27 | Discussions with my classmates through the LMS are less effective and less valuable |
| 28 | I can express myself freely through the activities conducted through the LMS |
| 29 | Learning through the LMS motivates me to learn more |
| 30 | I prefer to use and access learning content provided through the LMS than any other options |
| 31 | I am able to lean whenever and wherever I want once the learning content is made available through the LMS |
| 32 | Learning through the LMS reduces the interaction between students and lecturers inside classrooms |
| 33 | I will recommend to my classmates to study courses that are provided through the LMS |
| 34 | I believe that learning through LMS provides me with the opportunity to think critically and reach conclusions |
| 35 | Courses that are provided through the LMS are more interesting than other traditionally delivered courses |
| 36 | I believe that learning through LMS allows me to be more innovative |
| 37 | I feel that learning through the LMS requires advanced skills that I do not have |
| 38 | I believe that I will achieve great success if I continue studying all my courses through the LMS |
| 39 | Learning through the LMS increases the social bonding with my classmates |
| 40 | I am totally satisfied with my experience of learning some courses through the LMS |

# Implementation of Machine Learning Classification Algorithm Based on Ensemble Learning  for Detection of Vegetable Crops Disease

Pradeep Jha[1], Deepak Dembla[2], Widhi Dubey[3]
Department of Computer Science & Engineering, JECRC University, Jaipur, Rajasthan, India[1]
Department of Computer Application, JECRC University, Jaipur, Rajasthan, India[2]
Department of Botany, JECRC University, Jaipur, Rajasthan, India[3]

*Abstract*—In India, plant diseases pose a significant threat to food security, requiring precise detection and management protocols to minimize potential damage. Research introduces an innovative ensemble machine learning model for precise disease detection in tomato, potato, and bell pepper crops. Utilizing transfer learning, pre-trained models such as MobileNet and Inception are fine-tuned on a dataset of over 10,403 images of diseased and healthy plant leaves. The models are combined into a diverse ensemble, enhancing the precision and robustness of disease detection. The proposed ensemble models achieve an impressive accuracy rate of 98.95%, demonstrating their superiority over individual models in reducing misclassification and false positives. This advancement in plant disease detection provides valuable support to farmers and agricultural experts by enabling early disease identification and intervention.

*Keywords—DNN; transfer learning; crop; ensemble model; deep stacking and stacking approach; image pre-processing; tomato; bell paper; potato; disease*

## I. INTRODUCTION

In recent years, machine learning techniques have demonstrated significant potential in the field of plant disease detection [1]. This study introduces an innovative ensemble model that surpasses traditional approaches. This novel ensemble model combines several pre-trained models, including MobileNet, Inception, and ResNet, using a comprehensive dataset [2]. The primary objective of the model is to enhance the accuracy and resilience of conventional disease detection methods. This paper offers a thorough examination of the proposed model, elucidating its intricacies, and presents experimental results demonstrating its effectiveness in detecting plant diseases in tomato, bell pepper, and potato crops. In India, tomato, potato, and bell pepper are vital crops that make substantial contributions to the country's agricultural economy. However, these crops are vulnerable to various diseases that pose a significant threat to their yield. Notable tomato diseases in India include Spider mites, specifically the Two-spotted spider mite, Tomato mosaic virus, Target spot, Septoria leaf spot, Tomato Yellow Leaf Curl Virus, Late blight, Leaf Mold, Early blight, and Bacterial spot [3]. In contrast, potato crops are susceptible to fungal diseases like Early blight and Late blight, while bell pepper crops face the threat of bacterial spot disease. These collective diseases lead to substantial decreases in crop yield and quality, resulting

in economic hardships for farmers [4]. The Deep Stacking ensemble technique, as depicted in Fig. 1, involves training a set of base models and using their predictions as input for a higher-level advanced model. This advanced model, positioned at a higher level, is trained to merge the predictions derived from the base models, resulting in a final prediction. This approach has demonstrated its ability to enhance the accuracy and resilience of models. Nevertheless, the proposed ensemble model surpasses both of these methodologies in terms of accuracy.

Plant diseases pose a significant menace to global agricultural output and food security. The precise and early detection of plant diseases is of paramount importance for effective disease management and the mitigation of crop losses. This research paper introduces a comprehensive methodology for detecting diseases in three pivotal plant species: potato, tomato, and bell pepper. These crops play vital roles in global food production but are susceptible to a wide range of diseases that can significantly impact yield and quality. The proposed model leverages ensemble learning, deep stacking, and transfer learning, employing well-established convolutional neural network architectures, including MobileNe_v2, ResNet_v2, and Inception_v3 models. Ensemble learning enables the combination of multiple models, harnessing their diverse capabilities to enhance overall detection accuracy and robustness [5]. Deep stacking allows the capture of intricate interactions among the predictions of different models, thereby improving disease classification performance. Transfer learning plays a crucial role in the methodology. By leveraging pre-trained models trained on extensive image datasets like ImageNet, Proposed model can utilize their learned features and representations, which are transferable to various visual recognition tasks. Fine-tuning these models on plant disease dataset enables us to adapt their knowledge specifically to the challenges of disease detection in potato, tomato, and bell pepper plants. The objective of the model is to develop an automated and accurate disease detection system tailored to these specific plant species by harnessing the power of machine learning techniques, particularly ensemble learning and deep stacking. The significance of this work lies in its potential to elevate disease management practices within agriculture. By enabling early disease detection, farmers and agricultural experts can swiftly implement precise treatments and crop protection strategies,

thereby minimizing the adverse effects of diseases on both crop yield and quality.

Within this research paper, an elaborate exposition of the employed methodology is provided, encompassing the application of ensemble learning and deep stacking techniques using transfer learning models like MobileNet_v2, ResNet_v2, and Inception_v3. The assessment of performance is discussed, exploring practical implications of the proposed model's disease detection system in real-world agricultural contexts, and shedding light on prospective avenues for future research to enhance the system. Ultimately, this research aims to contribute to the advancement of automated disease detection systems in agriculture, with a specific focus on potato, tomato, and bell pepper plants. By harnessing the power of ensemble learning, deep stacking, and transfer learning, the aim is to improve disease management practices, increase crop productivity, and ensure food security in the face of plant diseases.



Fig. 1.   Block diagram of ensemble learning model.

## II.   RELATED WORK

First, Tulshan et al. [1], the authors presented a study on the detection of various plant leaf diseases. They focused on detecting diseases such as Early Blight, Mosaic Virus, Down Mildew, White Fly, and Leaf Miner. The authors achieved a high accuracy of 98.56% with their proposed K-Nearest Neighbors (KNN) model. This model was trained on a dataset consisting of 75 images. Furthermore, they compared the performance of their KNN model with a basic Support Vector Machine (SVM) model, which was trained on a larger dataset of 150 images. The SVM model achieved an accuracy of 97.6%.This research highlights the effectiveness of the KNN model in accurately detecting plant leaf diseases with a relatively smaller training dataset, outperforming the SVM model trained on a larger dataset. Ramesh et al. [6], the authors addressed the issue of detecting and classifying Rice Blast disease. Their approach utilized an Artificial Neural Network (ANN) algorithm. The training phase of their model achieved an accuracy of 99%, while the testing phase yielded a respectable accuracy of 90%.The authors gathered a dataset consisting of 300 images. These images were subjected to preprocessing steps, which involved conversion to the HSV color space and subsequent K-means clustering. It's worth noting that the dataset used in this research was their own, emphasizing the unique contribution of their work in the context of Rice Blast disease detection and classification. Lee et al. [7], conducted research on disease detection in potato

leaves, specifically targeting Potato Early Blight and Late Blight. They employed Convolutional Neural Networks (CNN) as their chosen algorithm, achieving an impressive accuracy rate of 99.09% in their experiments. The authors collected and utilized a dataset comprising 2150 images for their research. Furthermore, they conducted a comparative analysis by evaluating CNN against other machine learning algorithms such as Artificial Neural Networks (ANN), Backpropagation Neural Networks (BPNN), K-Nearest Neighbors (KNN), Support Vector Machines (SVM), among others. The results, presented in a table, demonstrated that the CNN model outperformed all other algorithms, yielding the highest accuracy of 99.09%.This survey provides valuable insights into the effectiveness of CNN in the detection of Potato Early Blight and Late Blight diseases in potato leaves, showcasing its superiority over various alternative methods.Asifet.al. [8],the authors explored various Convolutional Neural Network (CNN) architectures, including AlexNet, VggNet, ResNet, LeNet, and a Sequential model, for the detection of diseases on potato leaves. The specific diseases targeted were early blight, late blight, and Septoria blight. Their research achieved an accuracy rate of 97.00% in disease detection. To support this, they used a dataset containing 3000 images and employed a basic CNN approach, which involved image processing techniques and data augmentation to enhance training. It's worth noting that the training process for this model was conducted manually, underscoring the meticulous effort and expertise invested in the research. This study provides valuable insights into the application of various CNN architectures for effective disease detection on potato leaves, showcasing a high level of accuracy in identifying these plant diseases. Tembhurne et al. [9], the authors conducted research on the detection of plant diseases, specifically focusing on early blight and late blight. Their study utilized two deep learning models, namely AlexNet and GoogLeNet, as image classifiers. They evaluated the performance of these models on five key parameters: Accuracy, Precision, Sensitivity, F1 score, and Specificity. To support their research, the authors employed a dataset containing 900 images. The experimental outcomes showcased a remarkable 98.52% accuracy in identifying various plant diseases. Notably, the research was carried out on Kaggle, a widely recognized platform for data science and machine learning competitions and collaborative projects. This investigation offers valuable perspectives on the efficacy of deep learning models such as AlexNet and GoogLeNet in the realm of plant disease detection, highlighting their high accuracy and the use of comprehensive evaluation metrics. Naveen kumar et al. [10], the authors focused on the classification of plant diseases, specifically targeting early blight and late blight. For their research, they employed the InceptionResNetV2 model, achieving a commendable accuracy rate of 95.30%. To support their findings, the authors conducted a comparative analysis with other deep learning architectures, including VGG (16 and 19), ResNet 50, and InceptionV3. Interestingly, their results revealed that InceptionResNetV2 outperformed these models, showcasing a superior accuracy rate of 95.3%.It's worth noting that the dataset used in this research was proprietary, emphasizing the unique contribution of their study in the field of plant leaf disease classification. This research underscores the

effectiveness of InceptionResNetV2 as a robust model for accurate and reliable plant disease classification. Kukreja et al. [11], the authors addressed the challenge of potato disease classification, specifically focusing on Potato Early Blight and Late Blight, their approach employed a simple CNN based deep learning model. Through their research, they achieved a notable accuracy rate of 94.77%. The authors not only discussed basic classification results but also provided insights into classification metrics such as F1 score and recall. This study was conducted using a dataset consisting of 900 images and were part of the Plant Village project, emphasizing its contribution to the broader field of plant disease classification. The research highlights the effectiveness of their CNN-based model in accurately classifying potato diseases, particularly Early Blight and Late Blight. Tiwari et al. [12], the authors tackled the challenge of detecting potato leaf diseases, specifically focusing on Potato Early Blight and Late Blight. Their research introduced a novel approach by proposing a VGG19-based model. This model utilized VGG19 as a feature extractor, building upon the CNN architecture. In their experiments, the authors compared the performance of their VGG19-based model with VGG16 and Inception models. The dataset used for their study consisted of 2152 images, and their approach yielded an impressive accuracy rate of 97.80%. This research was conducted as part of the Plant Village project, emphasizing its contribution to the field of plant disease detection. The study underscores the effectiveness of their VGG19-based model in accurately detecting Potato Early Blight and Late Blight in potato leaves. Sinshaw et al. [13], the authors addressed the task of detecting Potato Late Blight disease. They employed three pre-trained CNN models, namely VGG16, VGG19, and InceptionV3, to explore the potential of transfer learning and data augmentation in disease detection. The results of their experiments demonstrated varying levels of accuracy among the three models. InceptionV3 emerged as the top-performing model, achieving an accuracy rate of 94.11%. In comparison, VGG16 achieved an accuracy of 93.2%, and VGG19 achieved 92.9%. These findings indicated that InceptionV3 outperformed VGG16 and VGG19 in detecting Potato Late Blight disease. The authors conducted their research using a dataset comprising 430 images, and it's noteworthy that the dataset was developed in-house, emphasizing the unique contribution of their study. This research highlights the effectiveness of transfer learning and data augmentation techniques, particularly when applied to the InceptionV3 model, in accurately detecting Potato Late Blight disease. Lakshmanarao et al. [14], the authors focused on predicting and classifying plant diseases, particularly within 15 classes encompassing Tomato, Pepper, and Potato. Their research employed ConvNets as the primary methodology. The study produced impressive accuracy rates for disease prediction and classification, with a notable 95% accuracy for Tomato, 98.5% for Bell Pepper, and 98.3% for Potato. To support their findings, the authors worked with a substantial dataset comprising 20,638 images. This research was conducted as part of the Plant Village project, underlining its valuable contribution to the field of plant disease prediction and classification. The study emphasizes the effectiveness of ConvNets in achieving accurate disease identification within the Tomato, Pepper, and Potato plant categories. Karthik et. al.

[15], the authors focused on predicting diseases in Tomato and Potato plants and their potential health benefits. Their research utilized a basic CNN approach, involving feature extraction and CNN model generation. This approach yielded a high accuracy rate of 98% in disease prediction for Tomato and Potato plants. The study was conducted as part of the Plant Village project, emphasizing its contribution to the field of plant disease prediction and its potential impact on the health benefits associated with these plants. The research underscores the effectiveness of their CNN-based deep learning techniques in accurately predicting diseases in Tomato and Potato plants.

Crop sustainability and food security pose pressing challenges in agricultural development. The early detection of plant diseases holds paramount importance in addressing these challenges. This study focuses on three crucial plant species: potato, tomato, and bell pepper, which are integral components of the global diet. It addresses the detection of seven different types of diseases. Plant diseases have garnered increasing attention in recent years, posing a significant threat to food security in India and around the world. Detecting these diseases demands expertise, time, and substantial human effort, as it requires proficiency in disease identification. Early disease detection is critical for minimizing damage and increasing crop yields, with far-reaching implications for the global economy. To tackle this issue, a novel ensemble model for plant disease detection is proposed in this research. This model leverages the capabilities of MobileNet, Inception, and ResNet models, which have been pre-trained on extensive image datasets, making them ideal for feature extraction in the proposed framework.

In this research, ensemble learning technique for vegetable leaves disease detection is implemented in the Google Colab environment, which utilizes GPU and TPU. The dataset comprises 10,403 images collected from primary sources and PlantVillage. Diverse image pre-processing methods, encompassing techniques like augmentation, Canny edge detection, noise reduction, and others, are utilized to extract distinctive features from images of plant leaves. Furthermore, implement transfer learning algorithm and then design Ensemble learning Algorithm. Their performance on various evaluation metrics, such as accuracy, f1-score, recall, and precision, is compared. These results are then compared with those obtained from the proposed ensemble model.

While the studies discussed have made significant contributions to plant disease detection, they also exhibit certain limitations that make them less suitable for our problem. For instance, the models proposed by Tulshan et al. [1], Ramesh et al. [6], and Lee et al. [7] were trained on relatively small datasets, which may limit their generalize ability to larger, more diverse datasets. The models used by Asif et al. [8] and Sinshaw et al. [13] demonstrated varying levels of accuracy, indicating a lack of consistency in their performance. Furthermore, the models employed by Naveenkumar et al. [10] and Kukreja et al. [11] were trained on proprietary datasets, which may limit the reproducibility of their results. Our proposed ensemble machine learning model addresses these limitations by leveraging the strengths of multiple pre-trained models and fine-tuning them on our extensive dataset. This approach enhances the precision and

robustness of disease detection, making it more suitable for the problem at hand.

### III. PROPOSED MODEL

The proposed model for this research is an ensemble machine learning model that integrates three deep neural network models: MobileNet, Inception, and ResNet. The choice of these models is motivated by their proven effectiveness in image classification tasks, their ability to handle large datasets, and their robustness to variations in disease presentation. By leveraging transfer learning, these pre-trained models are fine-tuned on our extensive dataset of diseased and healthy plant leaves, allowing them to capture intricate patterns and features specific to each disease. The ensemble approach strategically combines the strengths of individual models, enhancing the precision and robustness of disease detection. This methodology represents a significant advancement in the field of plant disease detection, providing valuable support to farmers and agricultural experts by enabling early disease identification and intervention.

#### A. Proposed Algorithm

Step 1: First load the acquired data using the imread() method present in CV2:

Step 2: Segment the data into training and testing Sets: from tensorflow.keras.preprocessing.image import ImageDataGenerator

train_datagen = ImageDataGenerator (rescale=1/255,shear_range=0.2, zoom_range=0.2,

horizontal_flip=True, validation_split=0.1)

Step 3: Normalize the data to remove anomalies

Step 4: Perform data augmentation to make the model robust:

To perform edge detection: cv2.Canny(img, 100, 200)

To flip the image: cv2.flip(img, 0)

To blur the image: cv2.blur(img, (20, 20))

For convolution   kernel = np.ones((7, 7), np.float3

conv = cv2.filter2D(img, -1, kernel)

Step 5: Implement Transfer learning, with models like MobileNet, Inception, ResNet

feature_extractor_model="https://tfhub.dev/google/imagen et/inception_v3/classification/5"pretrained_model_without_top _layer=hub.KerasLayer(feature_extractor_model,input_shape= (224, 224, 3), trainable=False)

similarly, used mobile net and ResNet models

Step 6: During this process evaluate all the models separately, using a confusion matrix, and compared their results using the respective recall, precision, and f1-scores.

Step 7: Combine the power of all the transfer learning models by using the ensemble learning technique used a deep stack library to combine the above models. Next, two classifiers, namely Random Forest Classifier and Extra Trees

Classifier, are specified as estimators for the second-level meta-learner. The meta-learner is the model responsible for amalgamating the predictions generated by the base models. It accepts the outputs or predictions from these base models as inputs and adapts its learning based on them how to best combine them to produce an improved prediction or decision and this is done using random forest.

Step 8: The evaluation measures to proposed model to compare it with the above models.

#### B. Data Collection

The foundation of the suggested plant disease detection model primarily relied on a dataset containing ten distinct classes. This dataset was derived from both Plant Village and a primary source. To facilitate for both the training and testing phases of the proposed model, a dataset was divided, allocating 80% for training and 20% for testing.

A thorough depiction of the dataset employed in this study is available in Table I.

TABLE I.    DETAILS OF DATA SET USED FOR PROPOSED MODEL TRAIN AND TEST PURPOSE

| S. No. | Crop Name | Disease | Test Image Data Set | Train Image Data Set |
|---|---|---|---|---|
| 1 | Pepper bell | Pepper bell Bacterial Spot | 201 | 697 |
| 2 | Pepper bell | Pepper bell Healthy | 297 | 1034 |
| 3 | Potato | Potato Late Blight | 60 | 300 |
| 4 | Potato | Potato Healthy | 31 | 106 |
| 5 | Potato | Potato Early Blight | 60 | 300 |
| 6 | Tomato | Tomato Spider Mites | 336 | 1173 |
| 7 | Tomato | Tomato Bacterial Spot | 437 | 1488 |
| 8 | Tomato | Tomato Leaf Mold | 191 | 666 |
| 9 | Tomato | Tomato Septoria Leaf Spot | 355 | 1239 |
| 10 | Tomato | Tomato Healthy | 319 | 1113 |

#### C. Pre-Processing

Data preprocessing is a pivotal phase in machine learning that readies raw data for analysis and model training. It heightens the quality and dependability of machine learning models by eliminating noise, rectifying errors, and converting data into a format that aligns with the model's requirements. In the context of colored images, data preprocessing is even more important due to the complex nature of image data.

- Noise removal: Colored images frequently exhibit noise, which can pose challenges for machine learning models in extracting features and achieving precise predictions. Data preprocessing techniques such as Gaussian blurring and median filtering can be used to remove noise from images, improving the performance of machine learning models.

- Color space conversion involves representing colored images in various color spaces, such as RGB, HSV, and YCbCr. Different color spaces have different properties, and some color spaces may be better suited for a particular machine-learning task than others. Data preprocessing techniques can be used to convert images from one color space to another, making it easier for machine learning models to learn from the data.

- Image resizing: Machine learning models often require images to be a certain size. Data preprocessing techniques can be used to resize images to the correct size, improving the performance of the model.

### D. Data Normalization

Effective data normalization is a pivotal step in data preprocessing, significantly improving the accuracy of models for crops image disease detection. This involves both scaling and standardizing the input data. For three-channel (RGB) images, as depicted in Fig. 2, the usual procedure entails calculating mean values for the RGB (Red, Green, and Blue) channels across the entire image dataset, often utilizing list comprehension. Notably, the red channel values demonstrate concentration towards lower values and a slight positive skew. Conversely, the green channel values appear more evenly distributed, featuring a prominent peak around 135, indicating a higher prevalence of green in these images. Lastly, the blue channel values exhibit the highest level of uniformity with minimal skew, but they display notable variation among images.



Fig. 2. RGB (Red, Green, and Blue) channels of leaves image.

### E. Image Processing

*1) First, resize all the images to 224x224 pixels:* Using cv2.resize() and creating numpy arrays, these numpy arrays can be employed for exploratory data analysis. When working with three-channel (red, blue, green) images in a dataset, data normalization becomes a crucial step in data pre-processing. To avoid potential issues with unnormalized data, normalization should be performed for each channel separately as depicted in Fig. 3. This approach reduces the impact of outliers or extreme values in the data, thereby enhancing the overall model performance.

*2) Canny Edge Detection as shown in Fig. 4* is a prominent tool in image pre-processing that helps define the boundary of objects in an image. This efficient technique reduces noise in an image by allowing us to focus on the objects inside the specified boundary [16].



Fig. 3. RGB (Red, Green and Blue) distribution channel value.



Fig. 4. Canny edge detection image.

*3) Further transformations are applied,* such as flippingas shown in Fig. 5, which involves reversing the rows and columns to achieve horizontal and vertical flipping. [16].

$$Image = A_{ijk}$$

$$\text{Horizontal flip:} A_{ijk} \rightarrow A_{i(n+1-j)k}$$

$$\text{Vertical flip:} A_{ijk} \rightarrow A_{(m+1-i)jk}$$

Although flipping maintains the same structure and features, a more diverse dataset can be created using this augmentation.

*4) Convolution is another augmentation technique* that performs a basic mathematical operation, in this a 2-D matrix window(kernel) moves across the length and breadth of the images shown in Fig. 6. It can be called a sunshine effect; this helps in building a robust and accurate model.



Fig. 5. Flip image.

Fig. 6. Convolved image.

*F. Proposed Ensemble Model*

To create a deep ensemble model using the Dirichlet Ensemble class from the deep stack library, the process begins with the importation of the Keras Member class from deep stack library's base module. This class is instrumental in creating ensemble members from Keras models. Subsequently, three members are generated for the ensemble, with each one aligned to a fine-tuned pre-trained model. Each member is equipped with its corresponding pre-trained model, and data generators are customized for both the training and validation datasets. This methodology simplifies the generation of ensemble members trained on identical data, sharing a common architecture while having distinct initial weights. Following the successful creation of these ensemble members, the Dirichlet Ensemble class is then employed to establish and train the ensemble model using the train() method.

To accomplish this task, the necessary modules for constructing a stacked ensemble model are imported. These include the StackEnsemble class from the deep stack library and the ensemble module from scikit-learn, which houses the base learners and the stacking classifier.

Specifically, import the following:

- sklearn: The scikit-learn module.
- StackEnsemble: The StackEnsemble class from the deep stack library for creating the stacked ensemble model.

- RandomForestClassifier: A base estimator from the sci-kit-learn ensemble module.
- ExtraTreesClassifier: Another base estimator from the sci-kit-learn ensemble module.
- StackingClassifier: A meta-estimator from the sci-kit-learn ensemble module for stacking the base learners.
- LogisticRegression: A logistic regression classifier for the meta-learner in the stacking classifier.

The mentioned modules and classes can be employed to construct a stacked ensemble model, which consolidates predictions from multiple base learners through a meta-learner. The base estimators encompass a diverse range of machine learning algorithms, including decision trees, support vector machines, and neural networks. In the realm of deep stacking, the base learners frequently consist of deep neural networks trained on the same dataset with varied architectures or hyper parameters.

As illustrated in Fig. 7, the meta-learner serves as a dedicated machine learning algorithm, aiming to learn the optimal combination of predictions from base learners to enhance overall model performance. In the context of deep stacking, the meta-learner often takes the form of a simple linear model, such as logistic regression or a neural network. This meta-learner is fed with predictions from the base learners as input and gains the ability to assign weights, ultimately producing the final prediction.

In the final stage of our ensemble model, a weighted average approach is employed. This approach assigns different weights to the predictions made by each base learner in the ensemble. The weights are determined by the meta-learner based on the performance of each base learner during training. This means that predictions from more accurate base learners are given more importance in the final prediction. The use of a weighted average ensures that our ensemble model leverages the strengths of each individual model, leading to a more robust and accurate prediction.



Fig. 7. Proposed model layer weighting average.

*G. Proposed Model Training*

All the member models are being compiled using the compile() function takes several arguments:

- Optimizer: specifies the optimizer algorithm used during training.

- Loss:designates the loss function employed during training, which quantifies the disparity between predicted and actual outcomes.

- Metrics: specifies the evaluation metric(s) that will be used to measure the model's performance during training.

All the member models are being fitted using the fit() function:

- Train generator refers to the input data generator responsible for producing batches of training data

- The batch size defines the quantity of samples utilized per gradient update throughout training, which has been set to 32.

- Validation data is the input data generator that generates batches of validation data.

- Validation steps specify the number of validation steps to evaluate the model after each epoch, which is 2

- Verbose is a flag that controls the level of logging during training. A value of 1 means that progress will be displayed during training.

- The model is trained using 10 epochs.

## IV. RESULTS AND DISCUSSIONS

*A. Model Training and Evaluation*

To access the efficacy of the proposed disease detection model in comparison to established counterparts, an evaluation was undertaken on a comprehensive dataset featuring images of diseased plants. This dataset encompassed ten distinct diseases across three plant species: potato, tomato, and bell pepper. The consistent outcomes of these experiments underscore the superior performance of the suggested model in disease detection accuracy when contrasted with other existing models. In direct comparison to individual transfer learning models, the ensemble approach coupled with deep stacking techniques exhibited remarkable performance, achieving higher precision, recall, and F1 scores. This heightened accuracy can be attributed to the synergistic effect of ensemble learning and deep stacking, enabling the model to capture and leverage diverse patterns and features associated with diseases in the targeted plants. Moreover, the proposed model showcased enhanced robustness, resulting in a diminished risk of misclassification and false positives when compared to other models. The effectiveness of the proposed model in early disease identification underscores its potential for real-world applications in supporting farmers and agricultural experts. By facilitating timely interventions and early-stage disease management, the proposed model has the capacity to contribute significantly to improved crop yields, minimized losses, and enhanced food security.

*1) MobileNet:* Based on Fig. 8, which illustrates Accuracy vs Epoch, it can be inferred that the model performs admirably shortly after 10 epochs. Notably, the accuracy of both Training and Validation sets exhibits minimal disparity, indicating robust model performance. The peak training accuracy stands at 97.4%, while the validation accuracy reaches 96.4%.

Fig. 9 provides insight into loss vs epochs. In this graph, both the training and validation losses undergo a sharp decrease after 10 epochs. The minimum loss values are achieved, with Training at 7.0% and Validation at 9.0%.



Fig. 8. Accuracy graph of mobilenet model.



Fig. 9. Loss vs. Epochs mobilenet.

*2) Inception:* The accuracy vs epochs graph in Fig. 10 the peak training accuracy: 91.03% and validation:90.2% accuracy, with a high jump in accuracy just after 10 epochs loss vs epochs which are shown in Fig. 11 the rapid fall in the losses of training and validation data and with the having losses of 25% and 28% for training and validation respectively.



Fig. 10. Accuracy graph of inception model.

Fig. 11. Loss vs. Epochs of inception model.

*3) ResNet:* In Fig. 12, the accuracy vs. epochs plot reveals that within the first 2 epochs, there is a substantial increase in accuracy, ultimately reaching 94.66% for training and 94.05% for validation by the 10th epoch.

Fig. 13 displays the losses vs. epochs plot, where it is evident that low losses are achieved after 10 epochs for both Training and Validation. Specifically, Training records a 15.2% loss, while Validation exhibits a 15.3% loss after 10 epochs.



Fig. 12. Accuracy graph of resnet model.



Fig. 13. Resnet Loss vs. Epoch graph of resnet model.

*4) Proposed model:* In Fig. 14, the accuracy vs. epochs plot illustrates that after 10 epochs, there is a significant increase in accuracy. By the 10th epoch, the accuracy reaches 94.66% for training and 94.05% for validation.

Fig. 15 depicts the plot of losses vs. epochs. As observed, low losses are achieved after 10 epochs for both Training and Validation. Specifically, Training records a 15.2% loss, while Validation exhibits a 15.3% loss after 10 epochs.



Fig. 14. Accuracy graph of proposed model.



Fig. 15. Loss vs. Epoch graph of proposed model.

*B. Performance Analysis*

The experiments were carried out using the Keras framework and Tensor Flow-GPU on Google Colab. To gauge the effectiveness of the proposed methodology, various metrics, including accuracy, precision, F-score, recall, and loss, were employed.

Accuracy, encompassing both positive and negative classes, was computed to provide insights into the model's proficiency in categorizing images accurately. Precision, recall, and F-score were utilized to offer additional insights into the model's capacity to accurately predict both positive and negative classes.

*1) CNN confusion matrix:* In Fig. 16 the confusion matrix there are few classes which have been mispredicted significantly, such as Tomato Spider mites has been predicted Tomato healthy 58 times, therefore impacting the accuracy and precision also there is some mispredictions between the species since pepper bell bacterial spot has been predicted as Tomato Spider mites 16 times and rest there were some misclassifications between the tomato species. The weighted average f1 score was 0.9.



Fig. 16. Confusion matrix representation of the CNN model.

*2) Confusion Matrix MobileNet:* In Fig. 17 the confusion matrix there are a few classes that have been miss predicted slightly, such as Tomato Septoria Leaf has been predicted Tomato Bacterial Spot 17 times, and pepper bell healthy has been predicted as pepper bell bacterial spot 12 times and rest of misclassifications were less than four times. Therefore, the mispredictions were significantly less than the CNN model indicating higher accuracy and hence its weighted average f1 score was 0.96.



Fig. 17. Confusion matrix representation of the mobilenet model.

*3) Confusion matrix inception model:* Inception Models in Fig. 18 the confusion matrix some high misclassifications in Tomato Septoria and Tomato Spider mite classes with misclassifications ranging from 3-54 times. And from the above model's confusion matrix it can be inferred that these two classes are most difficult to predict. The inception model has been predicting other classes accurately like the mobile net model but these two classes namely Tomato Septoria and Tomato Spider mite where mispredicted significantly thereby impacting the models weighted average f1 score which comes out to be 0.9.



Fig. 18. Confusion matrix representation of the inception model.

*4) Confusion matrix resnet model:* In Fig. 19 confusion matrix shows some high misclassifications in Tomato Bacterial Spot and Tomato Leaf mold classes, which is different from the above models as they have been misclassifying some other classes and have been classifying these classes correctly. But the rest of the classes are correctly predicted and the mispredictions are less than equal to four times. Therefore, the weighted F1 score is 0.94 of the inception models.



Fig. 19. Confusion matrix representation of the resnet model.

*5) Confusion matrix proposed ensemble model:* In Fig. 20 the proposed ensemble Model's confusion matrix showed much fewer misclassification than any of the above models. Only two classes Tomato Bacterial spot and Tomato Leaf Mold were wrongly predicted into some other tomato diseased classes and that too only 15 and 12 times respectively, However, apart from these instances, the majority of mispredictions were limited to two or fewer, underscoring the proposed model's superior performance in comparison to the aforementioned models across nearly all classes. The weight average F1 score, precision and recall are 0.97 greater than all the above models, showing its state-of-the-art performance.



Fig. 20. Confusion matrix representation of the proposed ensemble model.

*C. Comparative Analysis of Various Machine Learning Models with Proposed Model*

Table II displays the outcomes derived from various models, each utilizing 10,403 images and undergoing 10 epochs for both training and validation. The results suggest that the Stack Ensemble model outperforms other models in terms of accuracy. Fig. 21. Show Comparative analysis of Obtained Results from the MobileNet, Inception, ResNet Models with Proposed Ensemble Model and Fig. 21 shows comparative analysis of proposed work with related Model.

TABLE II.    OBTAINED RESULTS FROM THE DIFFERENT TECHNIQUES

| Model | Images set | Classes of disease | Epoch | Validation Accuracy |
|---|---|---|---|---|
| CNN | 10,403 | 10 | 10 | 90.05% |
| MobileNet | 10,403 | 10 | 10 | 95.90% |
| Inception | 10,403 | 10 | 10 | 89.91% |
| ResNet | 10,403 | 10 | 10 | 93.80% |
| Stack Ensemble | 10,403 | 10 | 10 | 97.86% |
| Dirichlet Ensemble | 10,403 | 10 | 10 | 97.80% |
| Proposed Model | 10,403 | 10 | 10 | 98.95% |



Fig. 21. Comparative analysis of obtained results from the different machine learning models with the proposed model.

## V. CONCLUSION

The timely and precise detection and diagnosis of plant health issues play a pivotal role in ensuring global food production. This study offers a comprehensive evaluation of a plant health assessment system specifically designed for potato, tomato, and bell pepper plants. Through extensive experimentation and comparative analysis against existing models, the study emphasizes the exceptional performance of the proposed model in accurately identifying and categorizing plant health issues. The model has undergone training on a dataset containing 10,403 images. By adopting a combined approach that incorporates ensemble learning and deep stacking techniques, the model consistently outperforms individual transfer learning models. The model attains an outstanding accuracy rate of 98.95% in accurately identifying plant diseases. This notable enhancement can be credited to the synergistic collaboration of ensemble learning and deep stacking, enabling the model to discern and capitalize on diverse patterns and features linked to plant health issues in the specific crops. Additionally, the proposed model showcases improved reliability by reducing the risks of misclassification and false positives when compared to alternative models. This increased reliability holds particular significance in real-world applications, where the accurate and timely recognition of plant diseases is of paramount importance for farmers and agricultural specialists. The efficacy showcased by the proposed model in promptly identifying diseases underscores its potential to offer practical solutions for farmers. By enabling timely interventions and the management of diseases in their early stages, the model has the capacity to significantly enhance crop yields, mitigate losses, and contribute to overall food security. This research underscores the substantial benefits of the proposed plant health assessment model and highlights its potential for real-world implementation. The integration of advanced techniques, improved accuracy, and enhanced reliability positions the model as a valuable tool in the field of agriculture, empowering farmers and experts to make well-informed decisions and effectively address plant health concerns.

## REFERENCES

[1] A. S. Tulshan and n. Raul, "Plant leaf disease detection using machine learning", IEEE 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT), pp. 1-6, 2019.

[2] P. Jha, D. Dembla and W. Dubey, "Crop Disease Detection and Classification Using Deep Learning-Based Classifier Algorithm", Emerging Trends in Expert Applications and Security. ICETEAS 2023. Springer Lecture Notes in Networks and Systems, vol 682, 2023.

[3] S. Kumar, S. Singh and V. Singh, "Tomato And Potato Leaf Disease Prediction With Health Benefits Using Deep Learning Techniques", IEEE 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT), 2021.

[4] Artzai Picon, Maximiliam Seitz, Aitor Alvarez-Gila, Patrick Mohnke, Amaia Ortiz-Barredo, JoneEchazarra, "Crop conditional convolutional neural networks for massive multi-crop plant disease classification over cell phone acquired images taken on real field conditions", Sci Direct Comput Electron Agric 167, pp. 1–10, 2019.

[5] P. Jha, D. Dembla, W. Dubey, "Comparative Analysis of Crop Diseases Detection Using Machine Learning Algorithm", IEEE 3rd International Conference on Artificial Intelligence and Smart Energy, 2023.

[6] S. Ramesh and D. Vydeki, "Rice blast disease detection and classification using machine learning algorithm.", IEEE 2018 8th International Conference on Computing, Communication and Networking Technologies (ICCCNT), 2018.

[7] T. -Y. Lee, J. -Y. Yu, Y. -C. Chang and J. -M. Yang, "Health Detection for Potato Leaf with Convolutional Neural Network", IEEE 2020 Indo – Taiwan 2nd International Conference on Computing, Analytics and Networks (Indo-Taiwan ICAN), pp. 289-293, 2020.

[8] M. K. R. Asif, M. A. Rahman and M. H. Hena, "CNN based Disease Detection Approach on Potato Leaves", IEEE 3rd International Conference on Intelligent Sustainable Systems (ICISS), pp. 428-432, 2020.

[9] Jitendra Tembhurne, Tarun Saxena, Tausif Diwan, "Identification of Plant Diseases Using Multi-Level Classification Deep Model", International Journal of Ambient Computing and Intelligence, Vol. 13, Issue. 1, pp. 1-21, 2022.

[10] M. Naveenkumar, S. Srithar, B. Rajesh Kumar, S. Alagumuthukrishnan and P. Baskaran, "InceptionResNetV2 for Plant Leaf Disease Classification," 2021 Fifth International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC), pp. 1161-1167, 2021.

[11] V. Kukreja, A. Baliyan, V. Salonki and R. K. Kaushal, "Potato Blight: Deep Learning Model for Binary and Multi-Classification," 2021 8th International Conference on Signal Processing and Integrated Networks (SPIN), pp. 967-672, 2021.

[12] D. Tiwari, M. Ashish, N. Gangwar, A. Sharma, S. Patel and S. Bhardwaj, "Potato Leaf Diseases Detection Using Deep Learning," 2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS), pp. 461-466, 2020.

[13] N. T. Sinshaw, B. G. Assefa and S. K. Mohapatra, "Transfer Learning and Data Augmentation Based CNN Model for Potato Late Blight Disease Detection," 2021 International Conference on Information and Communication Technology for Development for Africa (ICT4DA), pp. 30-35, 2021.

[14] A. Lakshmanarao, M. R. Babu and T. S. R. Kiran, "Plant Disease Prediction and classification using Deep Learning ConvNets," 2021 International Conference on Artificial Intelligence and Machine Vision (AIMV), pp. 1-6, 2021.

[15] K. Karthik, S. Rajaprakash, S. Nazeeb Ahmed, R. Perincheeri and C. R. Alexander, "Tomato And Potato Leaf Disease Prediction With Health Benefits Using Deep Learning Techniques," IEEE Fifth International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC), pp. 1-3, 2021.

[16] P. Jha, D. Dembla and W. Dubey, "Deep learning models for enhancing potato leaf disease prediction: Implementation of transfer learning based stacking ensemble model", Multimed Tools Application , 2023.

# Revolutionizing Software Project Development: A CNN-LSTM Hybrid Model for Effective Defect Prediction

Selvin Jose G[1], Dr. J Charles[2]

Department of Computer Science, Noorul Islam Centre for Higher Education, Kumaracoil, Tamil Nadu, India[1]
Department of Software Engineering, Noorul Islam Centre for Higher Education, Kumaracoil, Tamil Nadu, India[2]

*Abstract*—**Within the domain of software development, the practice of software defect prediction (SDP) holds a central and critical position, significantly contributing to the efficiency and ultimate success of projects. It embodies a proactive approach that harnesses data-driven techniques and analytics to preemptively identify potential defects or vulnerabilities within software systems, thereby enhancing overall quality and reliability while significantly impacting project timelines and resource allocation. The efficiency of software development projects hinges on their ability to adhere to deadlines, budget constraints, and deliver high-quality products. SDP contributes to these objectives through various means. This paper introduces a novel SDP model that harnesses the combined capabilities of Convolutional Neural Networks (CNNs) and Long Short Term Memory (LSTMs) unit. CNNs excel at extracting features from structured data, enabling them to discern patterns and dependencies within code repositories and change histories. LSTMs, conversely, excel in handling sequential data, which is pivotal for capturing the temporal aspects of software development and tracking the evolution of defects over time. The outcomes of the proposed CNN-LSTM hybrid model showcase its superior predictive performance. Simulation results affirm the substantial potential of this model to bolster the efficiency and reliability of software development processes. As technology advances and data-driven methodologies become increasingly prevalent in the software industry, the integration of such hybrid models presents a promising avenue for continually elevating software quality and ensuring the triumph of software projects. In summary, the utilization of this innovative SDP model offers a transformative approach to efficient software development, positioning it as a vital tool for project success and quality assurance.**

*Keywords—Data driven software development; proactive defect identification; software quality; predictive analytics; software defect prediction; artificial intelligence; long short term memory*

## I. Introduction

In the modern world, software plays an essential role in every aspect of our daily existence. It includes defense, automobile, healthcare, insurance, finance, banking, telecommunication, government administration sectors. In other words, the normal functionality of these sectors gets affected with the software failure. Technical and managerial issues are the two different issues normally emerge during the software development process. Thirty percent of project failures occur mainly due to technical issues and 70% are management issues [1]. Some of the problems related to managerial issues are insufficient risk management, customer buy-in, limited project resources and inaccurate project structure etc. However, low product delay, high expense, and schedule delay are the issues encountered during software program development. Prior to the analysis of software project risk, an efficient risk mitigation scheme should be developed by the program developer. Based on the accurate management of risk, the success of the project can be determined [2].

SDP is a crucial aspect of modern software development, aimed at improving the effectiveness and efficiency of software projects. In an era where software plays an increasingly pivotal role in our daily lives, organizations strive to deliver high-quality software products while minimizing the time and resources invested in debugging and maintenance. In this context, SDP emerges as an indispensable tool for the modern software development landscape, fostering both agility and the delivery of robust software products. Any fault, error, mistake, in a computer program, or a defect or bug in the software can cause unexpected or inaccurate results which are otherwise called a software defect. In order to enhance software quality, high-risk components must be detected as soon as possible [3].

Software defects can lead to an increase in both the cost and time required for delivering the expected end product. Also identification and rectification of defects is a highly waste of time and a costly software process [4]. One of the persistent challenges within the Software Development Lifecycle (SDLC) has been the ability to predict and identify defects during the initial phases of a project. In the current situation, development of a fault-free software which is highly reliable is a difficult task, as the problems for which software is developed is more complex and the domains that are involved are constantly increasing to constraints such as uncertainty and development processes that are complex [5].

At first, a collection of project data takes place from the software repositories. From the data, factors are calculated. The locations are predicted through models, which have a better potential for the defects contained. Ultimately, using prediction models, various measures are evaluated, such as precision, recall, and explanative power [6].

SDP using Deep Learning (DL) holds immense significance in the realm of software project development. In an era where software systems underpin virtually every aspect

of modern life, ensuring their reliability and robustness is paramount. Deep learning methodology, leveraging their capacity to analyzes extensive datasets and detect subtle patterns, present a compelling approach for preemptively recognizing and addressing potential flaws before they escalate into severe problems. By harnessing the power of DL, software development teams can enhance their efficiency, reduce maintenance costs, and deliver higher-quality products to users, ultimately contributing to the successful and sustainable advancement of software projects in an increasingly interconnected world. The major contribution of the proposed work includes:

- Software defect prediction utilizing CNN-LSTM hybrid model.

- To employ CNN in feature extraction from structured data, particularly in code repositories and change histories.

- To employ LSTM in handling sequential data, emphasizing their pivotal role in capturing the temporal aspects of software development and facilitating the tracking of defect evolution over time.

The paper is systematized as follows: Section II offers are view of the existing literature and identifies areas where further research is needed. Section III outlines the methodology in detail. In Section IV, the comprehensive results of the suggested approach are discussed. Finally, in Section V and Section VI, the paper concludes with a discussion and summary respectively.

## II. BACKGROUND

### A. Literature Review

Lei Qiao et al. [7] introduced an innovative methodology employing DL approaches for the anticipation of software system defects. This novel approach involves training a DL approach to predict the number of defects in software. Notably, when compared to widely adopted approaches such as Support Vector Regression, Feature-based Support Vector Regression, and Decision Tree Regression, the suggested method demonstrated a substantial enhancement in performance on established datasets. The improvement is notably reflected in a notable reduction in mean square error, ranging from 3% to 13%, and an augmentation in the squared correlation coefficient.

Pan et al. [8] introduced a range of CodeBERT models, specifically designed for SDP. The proposed research involved conducting empirical studies to assess the effectiveness of these approaches in cross-version and cross-project SDP scenarios. The findings demonstrated that leveraging pre-trained CodeBERT models led to enhanced prediction accuracy and time savings. Additionally, incorporating sentence-based and keyword-based prediction approaches further improved the effectiveness of pre-trained neural language frameworks in the context of SDP.

Geanderson Esteves et al. [9], delved into the realm of SDP models, harnessing the power of an efficiently implemented XGBoost variant, known as US-XGBoost. This endeavor

generated a multitude of random models, each meticulously assessed for the accuracy and interpretability. The key take away from the findings is that SDP is inherently project-specific. This means that the features constituting the most effective models can significantly differ from one project to another. Hence, comprehending the determinants behind model decisions becomes particularly vital.

Lakshmi Prabha and N. Shivakumar [10] introduced a novel hybrid model that addresses the challenge of classifying massive datasets accurately. The proposed approach combines feature reduction using Principal Component Analysis (PCA) with an overall probability application to minimize data loss during PCA processing. The approach further employed a neural network classification method for program bug detection. The simulation results demonstrated the model's impressive efficiency, achieving an outstanding 98.70 percent Area under the Curve (AUC) accuracy, marking a substantial advancement over existing models.

Hao Wang et al. [11] introduced GH-LSTMs, a novel DL framework for detecting potential code defects within software modules. GH-LSTMs leverage hierarchical LSTM architecture to simultaneously extract semantic and traditional features. A gated merge mechanism was employed to dynamically optimize the fusion of these features. Subsequently, a fully connected layer utilizes the combined features for within-project defect prediction. Remarkably, GH-LSTMs outperform existing methods in terms of F-measure, particularly in non-effort-aware cases.

Bilal Khan et al. [12] presented a comprehensive analysis of seven widely employed Machine Learning (ML) approaches applied to SDP. These approaches encompass SVM, J48, RF, MLP, RBF, HMM, and CDT. The evaluation of these methods utilized various performance metrics, including MAE, RAE, RMSE, RRSE, recall, and accuracy. The findings from the experiments revealed that NB and SVM exhibited superior performance in terms of minimizing MAE and RAE, respectively.

Shuo Feng et al. [13], delved into the robustness of SMOTE-based oversampling methods. This work not only probed the stability of these techniques but also introduced a set of novel and stable SMOTE-based oversampling strategies aimed at enhancing the reliability. These stable techniques minimize the inherent randomness in SMOTE by sequentially selecting defective instances, utilizing a distance-based approach for choosing neighbor instances, and ensuring an evenly distributed interpolation process. The proposed approach supported the findings with both mathematical proofs and empirical investigations across 26 datasets using four common classifiers. The simulation results demonstrated that the effectiveness of stable SMOTE-based oversampling approaches surpasses that of traditional SMOTE-based approaches in terms of stability and effectiveness.

Somya Goyal [14] introduced a pioneering Neighborhood-based Under-Sampling (N-US) algorithm to address the challenge of class imbalance. The study aims to showcase the efficacy of this N-US framework in enhancing accuracy for predicting defective modules. The experimental results revealed that the N-US approach successfully reduces the

dataset size by 17.29% and lowers the Imbalance Ratio (IR) by 19.73%. Consequently, it plays a vital role in augmenting classifier performance.

Cong Jin [15] introduced an innovative distance metric learning framework that leverages cost-sensitive learning (CSL) to mitigate the challenges posed by class-imbalanced datasets. This novel method, initially developed to address class imbalance, assigns distinct weights to individual training classes. Subsequently, this CSL-based distance metric learning is integrated into the large margin distribution machine (LDM) to take over the conventional kernel function. Empirical results indicated that these enhancements enable CS-ILDM to exhibit not only excellent predictive performance but also the lowest misprediction cost.

Kun Zhu et al. [16] introduced an innovative feature selection algorithm called EMWS, which optimally chooses a compact set of closely related features tailored to each software project. This approach effectively harnesses the local search capabilities of simulated annealing to augment the relatively weaker exploitation performance of the Whale Optimization Algorithm (WOA) while simultaneously capitalizing on WOA's strong global search abilities to enhance SA's exploration capabilities. A hybrid deep neural network model was also proposed. Empirical results substantiate that EMWS and WSHCKE consistently outperform various methods in various experiments.

Ruba Abu Khurma et al. [17] introduced the Island Model as an enhancement to the Binary Moth Flame Optimization (BMFO) algorithm for addressing the Feature Selection problem in the context of SDP. This innovative approach segments the moth population into multiple islands, facilitating feature sharing among them through migration. This technique serves to bolster solution diversity and govern algorithm convergence. The experiments involved assessing the performance of KNN, NB and SVM classifiers with and without FS, using BMFO-FS, and employing Is BMFO-FS. Notably, across three experiments, the SVM classifier consistently outperformed others, closely followed by the NB classifier.

Shuo Feng et al. [18] introduced a novel oversampling approach known as Complexity-based Oversampling Technique (COSTE). Instead of relying on inter-instance distances, COSTE assesses instance complexity to guide the selection of candidates for generating synthetic instances. The study evaluated COSTE's effectiveness against four other oversampling techniques using various classifiers, including, KNN, MLP, SVM and RF, across 23 imbalanced datasets. Remarkably, the simulation findings consistently demonstrated that COSTE outperformed the other methods across all performance metrics, highlighting its superior performance.

Shiqi Tang et al. [19] introduced TSboostDF, an innovative transfer-learning algorithm designed to address the complex problem of CPDP (Cross-Platform Domain Prediction). TSboostDF effectively combines the BLS sampling method, which considers the sample's weight, with transfer-learning techniques to mitigate the limitations commonly associated with conventional CPDP algorithms. This novel approach has been demonstrated to outperform other CPDP algorithms that rely on transfer-learning methods, highlighting its superior performance in resolving this challenging problem.

Liu Yang et al. [20] introduced an innovative hybrid algorithm that combines the strengths of SSA and PSO. This research involved a comprehensive analysis of the merits and limitations of swarm intelligence algorithms, aiming to devise strategies for enhancement. Notably, the empirical findings demonstrated that the hybrid approach integrating SSA and PSO, as presented in this work, significantly enhances the precision of software reliability model estimation and forecasting. Specifically, the proposed study focused on estimating and predicting software defects using the well-known G-O model. Furthermore, a fitness function was introduced, which is capable of effectively managing the parameter 'b' during initialization by leveraging the maximum likelihood formula.

An algorithm was presented by Nassif et al. [21] that aims to accomplish two significant tasks: learning and prediction. This approach has a high efficiency for other issues, such as software defect prediction, while being widely utilized in information retrieval. In this paper, two common output metrics namely bug density bug count were used as goal variables to compare various models. Additionally, it looked at how eight models with Grid Search optimization were affected by the use of imbalance learning and feature selection. The FPA scores of the bug density results have significantly improved with the usage of imbalance learning; however, the improvement in the bug count results has not been as great. Last but not least, applying feature selection with LTR decreased the bug density metric's FPA score but had no effect on the bug count findings.

### B. Research Gap

SDP models offer valuable insights and benefits, but they also come with several limitations. Some of the major limitations of existing SDP models are discussed below. These models heavily rely on historical data, which may be incomplete, inconsistent, or biased. Poor data quality can induce to inaccurate predictions. Software defect datasets often have imbalanced class distributions, with a small number of defective occurrences compared to non-defective ones. This can lead to model bias and lower predictive accuracy. Software systems, tools, and development practices evolve over time. Models trained on historical data may not effectively adapt to new technologies and practices. Creating relevant features from code repositories and other software data is a complex and manual process. Feature engineering can be time-consuming and error-prone. Complex ML models can over fit the training data, making them less generalizable to new projects or software environments. These models identify correlations but not necessarily causation. Identifying the root causes of defects often requires domain expertise and additional analysis. Models may not be transferable to different software domains or projects due to the unique characteristics of each project. Software is continuously evolving, and defects can emerge or be resolved after the training data was collected, making predictions less accurate. Bias in training data and predictions can lead to discrimination or unfair treatment in software development processes. Training and deploying sophisticated ML models can necessitate remarkable

computational resources and expertise, which may not be readily available to all development teams. Addressing these limitations requires careful consideration and often a combination of techniques, including data preprocessing, model selection, and ongoing monitoring and validation of the model's performance. Despite these challenges, SDP models have the potential to significantly improve software development processes when appropriately applied and maintained.

## III. MATERIALS AND METHODS

A CNN- LSTM based hybrid DL model is developed and analyzed for SDP for effective software project development. The detailed block schematic of the suggested work is illustrated in Fig. 1. The initial step of the work involves the dataset collection. It is followed by data preprocessing techniques. The preprocessed data is separate into training set and test set. The proposed hybrid model is trained and validated utilizing the training data and test data. Finally, the performance of the SDP model is analyzed.

### A. Dataset Description

The proposed system utilizes SDP dataset collected from Open ML, an online platform and repository for ML datasets. The dataset contains 31 features of 224 instances. In this paper, the binary classifier is developed to predict software defects based on 31 inputs.

### B. Units Data Preprocesing and Exploratory Data Analysis

Data preprocessing is a pivotal stage in data analysis, encompassing the tasks of cleansing, restructuring, and organizing raw data to make it appropriate for analysis or ML model training. The quality and efficacy of a learning model are substantially influenced by proper data preprocessing. Key techniques involved in this process include data cleaning, data transformation, handling missing values, addressing duplicate entries, and managing outliers. Among these techniques, the management of missing values stands out as a critical step. It involves handling data points that lack complete or relevant information. Various strategies can be employed for this purpose. One approach is to eliminate rows or columns with an excessive number of missing values, particularly if they do not significantly contribute to the analysis. An alternative method is imputation, which involves filling in the gaps with estimated or calculated values based on the data's distribution. For numerical data, this can involve mean, median, or mode imputation, ensuring that the dataset is more robust and suitable for analysis or modeling.

Exploratory Data Analysis (EDA) serves as a vital initial step in the data analysis process, where data analysts and scientists employ both visual and statistical methods to delve into a dataset. Its primary goal is to unveil patterns, relationships, anomalies, and insights within the data. EDA entails a range of techniques, including data visualization tools such as histograms, scatter plots, and box plots, in combination with summary statistics like mean, median, standard deviation, and more. This multifaceted approach allows for a comprehensive understanding of data distribution, the detection of outliers, an evaluation of data quality, and the development of an intuitive grasp of the dataset's underlying structure. In

practice, EDA plays a pivotal role in hypothesis formulation, guiding subsequent analytical processes, and informing decisions related to data preprocessing and modeling strategies. Ultimately, it aids in the discovery of valuable information and concealed patterns within the data. Summary statistics, which provide a concise summary of key dataset characteristics, include measures such as mean, median, mode, standard deviation, variance minimum and maximum values, quartiles including the first, second or median, and third quartiles, and counts or proportions for categorical variables. The summary statistics of SDP dataset is illustrated in Fig. 2.

EDA relies heavily on the strategic use of graphical representations to delve into a dataset. Through the construction of diverse charts, plots, and graphs, including histograms, scatter plots, box plots and heatmaps, it becomes possible to visually examine data distributions, reveal underlying patterns, identify anomalies, and grasp the interplay between variables. The data visualization of counts of classes in the dataset is visualized in Fig. 3.



Fig. 1. Block Schematics of Proposed SDP model.

| | id | total_loc | blank_loc | comment_loc | code_and_comment_loc | executable_loc | unique_operands | unique_operators | total_operands | total_operators | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 224.000000 | 224.000000 | 224.000000 | 224.000000 | 224.000000 | 224.000000 | 224.000000 | 224.000000 | 224.000000 | 224.000000 | ... |
| mean | 112.500000 | 26.013393 | 0.200893 | 5.607143 | 0.040179 | 20.205357 | 14.642858 | 10.473214 | 32.433037 | 49.834820 | ... |
| std | 64.807404 | 24.408239 | 0.862679 | 5.533197 | 0.218417 | 20.459360 | 10.056662 | 4.780303 | 29.730251 | 46.069885 | ... |
| min | 1.000000 | 2.000000 | 0.000000 | 0.000000 | 0.000000 | 2.000000 | 1.000000 | 2.000000 | 1.000000 | 3.000000 | ... |
| 25% | 56.750000 | 8.000000 | 0.000000 | 1.000000 | 0.000000 | 6.000000 | 8.000000 | 7.000000 | 14.000000 | 22.000000 | ... |
| 50% | 112.500000 | 19.000000 | 0.000000 | 5.000000 | 0.000000 | 12.000000 | 11.000000 | 9.000000 | 21.000000 | 31.000000 | ... |
| 75% | 168.250000 | 34.000000 | 0.000000 | 9.000000 | 0.000000 | 25.000000 | 18.000000 | 15.000000 | 50.000000 | 70.000000 | ... |
| max | 224.000000 | 95.000000 | 9.000000 | 30.000000 | 2.000000 | 82.000000 | 47.000000 | 19.000000 | 118.000000 | 184.000000 | ... |

8 rows × 31 columns

Fig. 2. Summary statics of SDP dataset.



Fig. 3. Data visualization of SDP dataset.

Fig. 4. Heatmap visualization of dataset.

Heatmap visualization serves as a powerful graphical tool for depicting data, representing information within a matrix by employing a range of colors. Its primary utility lies in the effective exploration of complex datasets, as each individual cell in the matrix corresponds to a specific data point, and the color intensity within these cells conveys the underlying data values, often using a color spectrum to illustrate the transition from lower to higher values. Fig. 4 illustrates the heatmap visualization of dataset.

*C. SDP Model using CNN-LSTM Hybrid Model*

The proposed defect prediction model is a hybrid system that leverages both CNNs and LSTM neural networks. This innovative approach aims to enhance the accuracy of identifying defects within software. CNNs are employed to detect spatial patterns in the code, such as relationships between different code segments, while LSTMs excel at modeling sequential dependencies over time. By amalgamating these two architectural components, the model becomes capable of understanding both localized and global patterns in the code base, enabling it to effectively pinpoint software defects. The approach involves encoding code as input sequences, which are then processed by CNN layers to capture spatial characteristics and subsequently by LSTM layers to capture temporal relationships. The resulting hybrid model offers a more precise and comprehensive defect prediction, which, in turn, supports the development of more dependable software systems.

A CNN is a type of deep ANN that emulates the human visual perception process, making it particularly effective for analyzing visual data [22]. CNNs leverage various multi-layer perceptron algorithms to reduce the need for extensive preprocessing of input data, making them widely adopted in the field of DL. CNNs represent one of the most prevalent neural network architectures, typically comprising millions of interconnected neurons organized into hierarchical layers, as illustrated in Fig. 5.



Fig. 5. Basic architecture of CNN.

CNN consist of three fundamental layers: convolutional, pooling, and fully-connected layers. The hidden layers within the convolutional and FC layers play a pivotal role in accessing the network's capacity to learn. The depth of a CNN is defined by the number of layers it comprises, and the deeper these layers are, the more intricate and abstract features they can extract from the input data, especially in the context of high-resolution images [23]. Within the CNN processing pipeline, the neurons in the input layer respond to visual stimuli, initiating the feature extraction process. The major aim of the convolutional layer is to capture these image features and propagate them to the subsequent hidden layers for computation, culminating in the extraction of results from the output layer. Activation functions often act as intermediaries between hidden layers, facilitating the transfer of valuable and essential information to inform the subsequent layers in the network.

The Long Short-Term Memory unit a prominent component of DL belongs to the family of Recurrent Neural Networks (RNNs) [24]. This specialized RNN excels in understanding and capturing intricate order dependencies within sequence prediction tasks. It is specifically designed to manage long-term relationships and tackle challenging problems, particularly those where the input order plays a pivotal role. Over time, numerous variations of Deep RNNs have been devised to combat issues related to vanishing and exploding gradients. Among these, the LSTM network stands out as a unique solution. It achieves its exceptional capabilities by employing distinct activation functions for each of its gates, allowing it to remember essential information from the past while efficiently discarding irrelevant data. Furthermore, LSTM incorporates an internal cell state vector, which serves as a practical representation of the network's retained knowledge from prior inputs. The LSTM unit employs three distinct gates: Forget Gate (f), Input Gate (i) and the Output Gate (o). Fig. 6 illustrates the core structural elements of an LSTM unit.



Fig. 6. Basic architecture of LSTM.

In a hybrid CNN-LSTM model, the initial CNN stage operates on the input sequence, effectively pinpointing spatial patterns or local features at each point in time. The resulting CNN layer outputs are subsequently fed into the LSTM component, which is responsible for modeling the temporal dependencies and capturing the sequential patterns within the data. This combined approach enables the model to effectively assimilate both local and global information, making it a robust choice for tasks that entail complex spatial and temporal interactions in sequential data. The architectural representation of the proposed model is visually depicted in Fig. 7.



Fig. 7. Proposed model architecture.

The suggested hybrid model begins with a dense layer featuring 100 neurons and a Relu activation function. This is succeeded by two additional dense layers, one with 50 neurons and the other with 25 neurons, both utilizing ReLU activation. The final dense layer has of a single neuron with a sigmoid activation function. Subsequently, an LSTM layer with 64 units is incorporated, followed by the application of a dropout layer with a 0.5 dropout rate to avoid over fitting. This is succeeded by another LSTM layer with 32 units. Finally, the model concludes with a dense layer consists a single neuron and a sigmoid activation function. This hybrid architecture combines a CNN with a recurrent neural network (RNN) using Time Distributed layers. The Time Distributed layers enable the application of feed forward network and flatten operations across each time step of the input sequence, effectively leveraging both spatial and temporal information for sequence-based tasks.

## IV. EXPERIMENTAL ANALYSIS

### A. Hardware and Software Setup

The proposed system utilizes SDP dataset contains 31 attributes. Google Collaboratory and Microsoft windows 10 are chosen in this research to ensure a stable computing experience. The system is equipped with an Intel Core i7-6850K 3.60 GHz 12- core processor, one NVIDIA Geforce GTX 1080 Ti GPU 2760 4MB. The dataset is split into two sets: training set (70%) and test set (30%). The entire procedure made use of Python and TensorFlow. The 'Adam' optimization function was used in the proposed model. The binary crossentropy is used as the loss function. The training process involved using a batch size of 32 for a total of 750 epochs. Finally, the proposed model predicts software defects as TRUE or FALSE.

### B. Result

The effectiveness of the suggested SDP model underwent an assessment using several key performance parameters, including accuracy, precision, recall, F1-score, specificity, and ROC AUC. Accuracy, a statistical measure, was employed to gauge the model's classification performance, representing the percentage of accurately predicted instances out of the entire dataset.

$$Accuracy = \frac{(TP+TN)}{(TP+TN+FP+FN)} \qquad (1)$$

Precision is a parameter utilized to find the model's capability to accurately predict positive outcomes. It quantifies the ratio of true positive predictions to all positive predictions, encompassing both accurate positive predictions and false positives.

$$Precision = \frac{(TP)}{(TP+FP)} \qquad (2)$$

Recall is a parameter that assesses a model's capacity to efficiently detect all pertinent examples of a specific class in a dataset. It quantifies the ratio of true positive predictions to the total number of actual occurrences belonging to that class.

$$Recall = \frac{(TP)}{(TP+FN)} \qquad (3)$$

The F1-Score serves as a metric employed to evaluate how well precision and recall are balanced in binary classification tasks.

$$F1-Score = 2 * \frac{(Precision*Recall)}{(Precision+Recall)} \qquad (4)$$

Specificity refers to the degree of precision or accuracy in targeting a particular characteristic, attribute, or aspect within a given context. The classification report of the proposed SDP model after simulation is tabulated in Table I.

An accuracy plot is a graphical representation that shows how well a model's predictions match the actual outcomes or labels in a dataset. It is a visual tool utilized to assess the effectiveness of a model, with the x-axis usually representing different iterations or epochs of training, and the y-axis indicating the accuracy of the model's predictions. The accuracy plot of the suggested prediction model is visualized in Fig. 8.

TABLE I. PERFORMANCE EVALUATION OF PROPOSED SDP MODEL

| Performance Metrics | Obtained Results |
|---|---|
| Accuracy | 95.58 % |
| Precision | 96.66 % |
| Recall | 93.54 % |
| F1- Score | 95.08 % |
| ROC AUC | 0.9542 |
| Specificity | 97.29 % |
| Cohens Kappa | 0.9108 |

Fig. 8. Accuracy plot of proposed SDP model.

Loss plot typically shows how the loss function, a measure of the error between the predicted and actual values, changes over time as the model learns from the training data. The loss plot is a critical tool for assessing the training process. The loss plot of the proposed SDP model is visualized in Fig. 9.



Fig. 9. Loss plot of proposed SDP model.

A confusion matrix is a visual representation of a classification model's performance that helps assess its accuracy and error rates. This grid is structured in such a way that it assigns actual class labels to its rows and predicted class labels to its columns. The cells along the diagonal of the matrix account for the accurate predictions, encompassing both true positives and true negatives. Meanwhile, any discrepancies outside the diagonal signify errors, including false positives and false negatives. The confusion matrix obtained for the proposed SDP model is illustrated in Fig. 10.

The Receiver Operating Characteristic (ROC) curve is a visual tool utilized to assess the effectiveness of models. It provides a graphical representation of how well a model can differentiate between positive and negative classes by depicting the balance between its true positive rate and false positive rate

across various threshold settings. When examining a ROC curve, a flawless model would closely follow the upper-left corner of the plot, demonstrating high sensitivity and minimal false positives, whereas random guessing would produce a diagonal line running from the lower-left to the upper-right. The Area Under the ROC Curve (AUC-ROC) serves as a succinct metric summarizing the model's overall performance, with greater values indicating superior discriminatory capabilities. The ROC curve is visualized in Fig. 11. Table II shows the performance comparison of proposed model with existing model.



Fig. 10. Confusion matrix of proposed SDP model.



Fig. 11. ROC curve.

TABLE II.    PERFORMANCE COMPARISON OF PROPOSED MODEL WITH EXISTING MODELS

| Author &Reference | Methodology | Result |
|---|---|---|
| 12 | RF | 88.32% |
| 14 | KNN | 94.6 |
| 16 | CNN and kernel extreme learning machine | 93.5 |
| 17 | moth flame optimization | 80% |
| 20 | Hybrid Particle Swarm Optimization and Sparrow Search Algorithm | 88% |
| **Proposed model** | **CNN+LSTM** | **95.58%** |

## V.    DISCUSSION

In the realm of software defect prediction, the findings of this study contribute valuable insights into the identification and mitigation of software defects. Through a comprehensive analysis and employing advanced learning algorithms, the study successfully establishes robust models for predicting potential defects in software development. The results reveal key predictors and patterns associated with defect occurrence, shedding light on critical areas that warrant attention during the software development life cycle. Feature extraction from structured data making them adept at identifying patterns and dependencies within code repositories, change histories. On the other hand, are well-suited for handling sequential data, which is crucial in capturing temporal aspects of software development and tracking the evolution of defects over time.

Moreover, the study's exploration of different feature sets and model evaluation techniques enhances the reliability of the proposed defect prediction models. The simulation results demonstrated that the adoption of a CNN-LSTM hybrid model for SDP has the potential to significantly contribute to more efficient and reliable software development processes.

## VI.    CONCLUSION

SDP plays a pivotal role in ensuring efficient development in software projects. This approach is proactive in nature, utilizing data-driven methods and analytics to detect possible flaws or weaknesses in software systems prior to the escalation into critical problems. This contributes not only to the improvement of the software's overall quality and dependability but also exerts a substantial influence on project schedules and resource allocation. Efficient software development projects are characterized by their ability to meet deadlines, stay within budget constraints, and deliver a high-quality product. SDP contributes to these goals in several key ways. This paper presented a SDP model utilizing both the benefits of CNN and LSTM. This approach leverages the power of CNN and LSTM to address the challenges associated with identifying and mitigating software defects, ultimately contributing to the improvement of software quality and project timelines. The CNN-LSTM hybrid model combines the strengths of both convolutional and recurrent neural networks. CNNs excel in feature extraction from structured data, making them adept at identifying patterns and dependencies within code repositories and change histories. LSTMs, on the other hand, are well-suited for handling sequential data, which is crucial in capturing temporal aspects of software development and tracking the evolution of defects over time. The proposed prediction model achieved better prediction results. The simulation results demonstrated that the adoption of a CNN-LSTM hybrid model for SDP has the potential to significantly contribute to more efficient and reliable software development processes. As technology continues to advance and data-driven approaches become increasingly prevalent in the software industry, the integration of such models holds promise for continually enhancing software quality and the success of software projects.

## REFERENCES

[1]   Verner, J., Sampson, J., & Cerpa, N. (2008, June). What factors lead to software project failure?. In 2008 second international conference on research challenges in information science (pp. 71-80). IEEE.

[2]   Sangaiah, AK, Samuel, OW, Li, X, Abdel-Basset, M & Wang, H 2018, 'Towards an efficient risk assessment in software projects-Fuzzy reinforcement paradigm', Computers & Electrical Engineering, vol. 71, pp. 833-846. https://doi.org/10.1016/j.compeleceng.2017.07.022.

[3]   Shukla, HS & Verma, DK 2015, 'A Review on Software DefectPrediction', International Journal of Advanced Research in ComputerEngineering & Technology (IJARCET), vol. 4 no. 12, pp. 4387-4394.

[4]   Gupta, V, Ganeshan, N & Singhal, TK 2015, 'Developing SoftwareBug Prediction Models Using Various Software Metrics as the BugIndicators', International Journal of Advanced Computer Science &Applications, vol. 1, no. 6, pp. 60-65.

[5]   Vashisht, V, Lal, M, Sureshchandar, GS & Kamya, S 2015, 'Aframework for software defect prediction using neural networks',Journal of Software Engineering and Applications, vol. 8, no. 8,pp. 384-394.

[6]   Shihab, E 2012, An exploration of challenges limiting pragmaticsoftware defect prediction (Doctoral dissertation, Queen's University).Link:http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.467.1447&rep=rep1&type=pdf.

[7]   Qiao, L., Li, X., Umer, Q., & Guo, P. (2020). Deep learning-based software defect prediction. Neurocomputing, 385, 100-110.

[8]   Pan, C., Lu, M., & Xu, B. (2021). An empirical study on software defect prediction using codebert model. Applied Sciences, 11(11), 4793.

[9]   Esteves, G., Figueiredo, E., Veloso, A., Viggiato, M., & Ziviani, N. (2020). Understanding machine learning software defect predictions. Automated Software Engineering, 27(3-4), 369-392.

[10]  Prabha, C. L., & Shivakumar, N. (2020, June). Software defect prediction using machine learning techniques. In 2020 4th International Conference on Trends in Electronics and Informatics (ICOEI)(48184) (pp. 728-733). IEEE.

[11]  Wang, H., Zhuang, W., & Zhang, X. (2021). Software defect prediction based on gated hierarchical LSTMs. IEEE Transactions on Reliability, 70(2), 711-727.

[12]  Khan, B., Naseem, R., Shah, M. A., Wakil, K., Khan, A., Uddin, M. I., & Mahmoud, M. (2021). Software defect prediction for healthcare big data: an empirical evaluation of machine learning techniques. Journal of Healthcare Engineering, 2021.

[13]  Feng, S., Keung, J., Yu, X., Xiao, Y., & Zhang, M. (2021). Investigation on the stability of SMOTE-based oversampling techniques in software defect prediction. Information and Software Technology, 139, 106662.

[14]  Goyal, S. (2022). Handling class-imbalance with KNN (neighborhood) under-sampling for software defect prediction. Artificial Intelligence Review, 55(3), 2023-2064.

[15]  Jin, C. (2021). Software defect prediction model based on distance metric learning. Soft Computing, 25, 447-461.

[16]  Zhu, K., Ying, S., Zhang, N., & Zhu, D. (2021). Software defect prediction based on enhanced metaheuristic feature selection optimization and a hybrid deep neural network. Journal of Systems and Software, 180, 111026.

[17]  Khurma, R. A., Alsawalqah, H., Aljarah, I., Elaziz, M. A., & Damaševičius, R. (2021). An enhanced evolutionary software defect prediction method using island moth flame optimization. Mathematics, 9(15), 1722.

[18]  Feng, S., Keung, J., Yu, X., Xiao, Y., Bennin, K. E., Kabir, M. A., & Zhang, M. (2021). COSTE: Complexity-based OverSampling TEchnique to alleviate the class imbalance problem in software defect prediction. Information and Software Technology, 129, 106432.

[19]  Tang, S., Huang, S., Zheng, C., Liu, E., Zong, C., & Ding, Y. (2021). A novel cross-project software defect prediction algorithm based on transfer learning. Tsinghua Science and Technology, 27(1), 41-57.

[20] Yang, L., Li, Z., Wang, D., Miao, H., & Wang, Z. (2021). Software defects prediction based on hybrid particle swarm optimization and sparrow search algorithm. Ieee Access, 9, 60865-60879.

[21] Nassif, A. B., Talib, M. A., Azzeh, M., Alzaabi, S., Khanfar, R., Kharsa, R., & Angelis, L. (2023). Software defect prediction using learning to rank approach. *Scientific Reports*, *13*(1), 18885.

[22] O'Shea, K., & Nash, R. (2015). An introduction to convolutional neural networks. arXiv preprint arXiv:1511.08458.

[23] Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," Nature, vol. 521, no. 7553, pp. 436–444, 2015.

[24] Sherstinsky, A. (2020). Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. Physica D: Nonlinear Phenomena, 404, 132306.

# US Road Sign Detection and Visibility Estimation using Artificial Intelligence Techniques

Jafar AbuKhait

Dept. of Computer and Communications Engineering, Tafila Technical University, Tafila, Jordan

*Abstract*—**This paper presents a fully-automated system for detecting road signs in the United States and assess their visibility during daytime from the perspective of the driver using images captured by an in-vehicle camera. The system deploys YOLOv8 to build a multi-label detection model and then, calculates various readability and detectability factors, including the simplicity of the surroundings, potential obstructions, and the angle at which the road sign is positioned, to determine the overall visibility of the sign. This proposed system can be integrated into Driver Assistance Systems (DAS) to manage the information delivered to drivers, as an excess of information could potentially distract them. Road signs are categorized based on their visibility levels, allowing Driver Assistance Systems to caution drivers about signs that may have lower visibility but are of significant importance. The system comprises four main stages: 1) identifying road signs using YOLOv8; 2) segmenting the surrounding areas; 3) measuring visibility parameters; and 4) determining visibility levels through fuzzy logic inference system. This paper introduces a visibility estimation system for road signs specifically tailored to the United States. Experimental results showcase the system's effectiveness. The visibility levels generated by the proposed system were subjectively compared to decisions made by human experts, revealing a substantial agreement between the two approaches.**

*Keywords—Road sign detection; YOLOv8; driver assistance system; fuzzy logic; detectability; visibility estimation*

## I. INTRODUCTION

In recent times, there has been a notable rise in the adoption of Driver Assistance Systems (DAS), primarily driven by the expanding complexity of road networks [1]. These systems are integrated into vehicles to simplify the driving experience and enhance overall driver safety. Road signs serve as a vital source of information for both drivers and these advanced systems, yet their visibility and detectability (the driver's capacity to spot a road sign within a complex or cluttered environment, essentially measuring how effectively the sign stands out) can be compromised in various scenarios. These scenarios can be categorized as either temporary, influenced by factors like lighting and adverse weather conditions, or permanent, resulting from vandalism or improper sign placement [2].

Reduced visibility of road signs significantly diminishes the effectiveness of communication between drivers and these signs. Consequently, DASs can play a crucial role in notifying drivers about warnings in such situations. Road sign detection is a basic step that every DAS system should have. It is noteworthy that an effective Driver Assistance System (DAS) should achieve a balance, avoiding the inundation of drivers

with excessive road information. Overloading drivers with information may pose a risk of distraction, as discussed in [3].

Employing computer vision techniques in Driver Assistance Systems (DASs) enables the detection and estimation of road sign visibility. This information can then be used to alert drivers about crucial warnings regarding less visible signs. The implementation of these techniques contributes to enhanced driver safety.

In this work, we propose a fully-automated computer vision system for detecting and assessing the visibility of road signs in the United States in terms of their detectability and readability. Detectability is defined as the driver's capacity to identify and acknowledge the presence of specific road signs within complex or cluttered environments while readability represents the clearness degree of the foreground text on the sign. This proposed system can be integrated into Driver Assistance Systems (DAS) to streamline the information presented to drivers. Furthermore, transportation agencies could leverage this system to assess the placement of road signs across their road networks. The proposed system deploys YOLOv8 in the detection of road signs and estimates their visibility using fuzzy logic after measuring five different visibility parameters.

The proposed system aims to: 1) implement a fully-automated multi-label detection model of United States road signs using YOLOv8 which is the latest YOLOs' detection algorithm; 2) measuring five novel visibility parameters of road signs that describe both sign readability and detectability; 3) evaluate visibility level of road signs to low, medium and high using fuzzy inference system that connects the suggested visibility parameters to the visibility output.

The rest of the article is prepared as follows: Section II provides a background of road sign detection and visibility estimation systems. Section III demonstrates the proposed visibility estimation system. Section IV evaluates the system performance experimentally based on certain metrics. Lastly, Section V elaborates the conclusions.

## II. RELATED WORK

Automated estimation systems for assessing the visibility of road signs should integrate Road Sign Detection (RSD) systems. The primary objective of RSD is to pinpoint the location of road sign objects within a scene or within images captured from inside a vehicle. RSD systems can be primarily categorized into two groups: those reliant on color and those based on shape. In the realm of color-based detection, some researchers have utilized RGB color thresholding to segment

road sign images, as demonstrated in [4, 5], while others have proposed the use of HSI color space for the same purpose, as indicated in [6].

Conversely, shape-based approaches have also been put forth by various researchers. In study [4], for instance, a Support Vector Machine (SVM) was trained using four vectors representing distances from the border to the bounding box to recognize road sign shapes. Researchers in [7] employed a Distance to Border (DtB) vector to identify the shape of road signs. For detecting Regions of Interest (ROI), a boosted detector cascade was trained using dissociated dipoles, while the recognition of triangular or circular road sign shapes was achieved through the utilization of the Hough transform and radial symmetry, as described in [8]. In study [9], a genetic algorithm was employed, and Haar-like features were deployed in study [6] to detect road sign shapes. Researchers in study [10] employed a set of cascaded geometric detectors, capitalizing on the inherent symmetry of road sign shapes for detection and recognition. In study [11], a speed recognition system has been proposed based on independent component analysis. In research [12], geometric features were deployed in the recognition of speed signs in the United States.

Recently, Convolutional Neural Network (CNN) was deployed with different architectures in detection tasks [13, 14] including road sign detection and recognition. In research [15], mask R-CNN was deployed to detect 200 traffic-sign categories with automatic end-to-end learning. In study [16], the authors analyzed seven architectures for detecting the road signs: YOLO, YOLOv2, YOLOv3, PP-YOLO model and R-CNN, Fast R-CNN, Faster R-CNN. This study implies that YOLOv3 and Faster R-CNN perform comparatively better for road sign detection. In [17], authors proposed a detection model to detect and identify traffic signs based on YOLOv7 and Convolutional Block Attention Module. In study [1], authors proposed a road sign detection and recognition system based onYOLOv5s object detection algorithm.

Several researchers have also proposed methods for estimating road sign visibility from digital images. Researchers in [18] introduced a novel technique for measuring road sign retro-reflectivity using two varied illumination images. In study [3], traffic signals' detectability and discriminability were quantified using in-vehicle images. Researchers in [19] utilized five image features to gauge the visibility of certain sign. For visibility estimation in foggy conditions, authors of study [20] introduced a method utilizing in-vehicle images. Researchers in [21] extracted both local and global features to evaluate a human driver's ability to detect and recognize road sign objects. Lastly, researchers in [22] showcased a novel approach based on SVM learning to estimate road sign saliency. In study [23], tilt angle of road sign was used to assess its condition. In study [24], various detectability features of road signs were measured to estimate the visibility level in cluttered environments.

In research [25], a three-dimensional approach for visibility assessment of highway signs has been proposed. The proposed approach measures sign's visibility, legibility, and readability based on its placement, height, and traffic flow. In [26], a system for classifying horizontal road signs as correct or with poor visibility is proposed. This system deploys YOLOv4-Tiny neural network model for classification and the contrast difference for visibility estimation. In [27], a study authors proposed a camera-based visibility estimation method for a traffic sign. The proposed method integrates both the local features and global features in a driving environment. These features measure sign's positional relationships and the contrast between a traffic sign and its surroundings. In research [28], author proposed an imaging-based system to estimate road sign visibility in a cluttered environment from the driver's perspective in daytime using in-vehicle camera images. The proposed system deploys a geometric sign detector and suggests two visibility parameters which are color difference and environment complexity. In study [29], authors proposed a method that can automatically detect the occlusion and continuously quantitative estimate the visibility of traffic sign. The proposed method deploys road sign orientation and occlusion in evaluating its visibility. In study [30], authors proposed a quantitative visual recognizability evaluation method for traffic signs in large-scale traffic environment. The proposed method evaluates the geometric, occlusion and sight line deviation factors of traffic signs.

In conclusion, the literature demonstrates that the implementation of automated vision-based road signs detection and recognition systems represents a significant advancement in modern transportation networks. In addition, visibility of these road signs is a major concern for both drivers and transportation agencies. The literature shows a lot of shortcomings of current road sign visibility estimation systems which can be concluded as the lack of automation in both the detection and visibility estimation, the deficiency of road sign detection models under different illumination and occlusion conditions, the failure to measure all visibility parameters that represent the road sign readability and detectability, and the need to estimate the road sign visibility to various levels either by a rule-based or machine learning techniques. As artificial intelligence techniques continue to evolve towards greater efficiency, these systems could be improved and automated completely for better safety over transportation networks. Additionally, current road sign visibility estimation systems should deploy powerful detection models that have the capability

## III. THE PROPOSED SYSTEM

The proposed system based on road sign imaging, depicted in Fig. 1, comprises four distinct modules:

*1) Multi-label road sign detection:* In this initial module, the system builds a detection model using YOLOv8 algorithm based on the in-vehicle images to detect and identify three categories of road sign objects (regulatory, warning and stop signs).

*2) Cropping of surrounding regions:* In this module, the system geometrically extracts four adjacent regions around the road sign object. These regions would be used in the next module to calculate some visibility parameters.

*3) Measurement of visibility parameters:* During this module, the system establishes and computes five visibility

parameters which are: 1) readability of sign foreground; 2) color difference between the sign and its four surrounding regions; 3) complexity of surroundings; 4) occlusion; and 5) tilting. These parameters characterize both the readability and detectability of each road sign.

*4) Determination of visibility levels:* In this module, the system assesses and categorizes the visibility level of each road sign as low, medium, or high using fuzzy logic inference system.

### B. Multi-Label Road Sign Detection

In this module, YOLOv8 algorithm is used to detect three different types of US road signs: 1) Regulatory Signs (White Rectangular-shape signs); 2) Warning Signs (Yellow Diamond-shape signs); and 3) Stop Sign (Red Octagonal-shape signs. The multi-label detection model was trained on Google Colab notebook after building an annotated road sign dataset of 664 images. Once the model was trained, it was tested on a separate set of validation images to evaluate its performance.



Fig. 1. Flow diagram of the proposed system.

*1) Dataset preparation:* Originally, 664 images in which one US road sign is existed, was collected and uploaded to Roboflow. Data analysis operations were achieved such as: pre-processing, resizing, annotation, augmentation and health check. All images were resized to 640x640. The following augmentations were done on these images: Rotation (between -21° and +21°), Saturation (between -20% and +20%), Brightness (between -20% and +20%), Blur (up to 1px), Noise (up to 3% of pixels). A set of 1519 images was achieved splitted as: 1332 for training, 116 for validation, and 71 for testing.

*2) Model training and evaluation:* The model was trained using YOLOv8. It was evaluated for detecting three classes: Stop signs, Warning signs and Regulatory signs. The number of Epochs used to train the model was 150. The model detection performance was evaluated using mean average precision (mAP), recall and precision.

The output of this module is road sign surrounded by a bounding box as shown in Fig. 2. This detected road sign would be used in the next modules to estimate its visibility.



Fig. 2. Examples of detected road signs.

### C. Cropping of Surrounding Regions

In the module of the proposed system, road sign visibility is characterized by the driver's capacity to distinguish the sign's location from the surrounding background in a real-life scenario. Various elements in the background might divert the driver's attention away from identifying the road sign's location. To gauge visibility, we assess the road sign's location in relation to its surroundings. For Stop, Warning, and Regulatory signs, we have extracted four adjacent regions from the input image, as illustrated in Fig. 3. This process has been accomplished by mirroring the bounding box in the four directions. For Warning signs, the four regions were obtained after rotating the sign. Each region possesses a symmetrical shape and double the area of the sign region. These four surrounding regions are denoted as R1, R2, R3, and R4, while the sign region is designated as S, as shown in Fig. 3.



Fig. 3. The four surrounding regions for: a) Regulatory sign; b) Warning sign; c) Stop sign.

## D. Measurement of Visibility Parameters

Different detectability and readability parameters of road sign region and surrounding regions are used to determine the visibility level of the road sign. Five parameters are proposed to describe the visibility of road signs: 1) readability of sign foreground; 2) color difference between the sign and the four surrounding regions; 3) complexity of surroundings; 4) sign occlusion; and 5) sign tilting.

Fig. 4 shows some road signs exhibiting poor visibility based on these parameters. Each parameter is designed as low or high based on ranges that were determined by a human expert.

*1) Readability of sign foreground:* This parameter measures the clearness degree of the foreground text on the sign. The greater the difference between the foreground and background is better considering the different colors in the three sign classes. This parameter, denoted by R, has been computed on the gray images of the detected signs by subtracting the average gray levels of both foreground and background as:

$$R = G_F - G_B \tag{1}$$

where, $G_F$ is the gray level of the sign foreground (white on Stop signs and black on Warning and Regulatory signs), $G_B$ is the gray level of the sign background (red on Stop signs, yellow on Warning signs and white on Regulatory signs).

A lower color difference between the sign foreground and background (0 - 120) diminishes a driver's ability to read and recognize road signs, whereas a higher color difference (120 - 255) enhances readability probabilities. Therefore, a significant color difference between the sign foreground and background leads to improved road sign visibility.



(a)   (b)   (c)

(d)   (e)

Fig. 4.    Examples of low visibility road signs due to: a) Color difference between sign and surroundings, b) Occlusion, c) Complexity of surrounding regions, d) Readability of sign foreground, e) Tilting.

*2) Color difference:* This parameter measures the clearness degree of the sign with respect to its surrounding regions. The process involves computing the average color of the RGB values for both the road sign and its four surrounding regions. The color disparity between the sign region and each of its four surrounding regions is then quantified as follows:

$$D1 = \sqrt{(\overline{R}_S - \overline{R}_{R1})^2 + (\overline{G}_S - \overline{G}_{R1})^2 + (\overline{B}_S - \overline{B}_{R1})^2} \tag{2}$$

$$D2 = \sqrt{(\overline{R}_S - \overline{R}_{R2})^2 + (\overline{G}_S - \overline{G}_{R2})^2 + (\overline{B}_S - \overline{B}_{R2})^2} \tag{3}$$

$$D3 = \sqrt{(\overline{R}_S - \overline{R}_{R3})^2 + (\overline{G}_S - \overline{G}_{R3})^2 + (\overline{B}_S - \overline{B}_{R3})^2} \tag{4}$$

$$D4 = \sqrt{(\overline{R}_S - \overline{R}_{R4})^2 + (\overline{G}_S - \overline{G}_{R4})^2 + (\overline{B}_S - \overline{B}_{R4})^2} \tag{5}$$

where, $(\overline{R}_S, \overline{G}_S, \overline{B}_S)$ are the average RGB colors of the sign region and $(\overline{R}_{Ri}, \overline{G}_{Ri}, \overline{B}_{Ri})$ are the average RGB colors of each surrounding region Ri.

Subsequently, these four disparity values are averaged to derive the overall color difference value, denoted as D. A lower color difference (0 - 120) diminishes a driver's ability to detect road signs, whereas a higher color difference (120 - 255) enhances detection probabilities. Therefore, a significant color difference between the sign region and its adjacent regions leads to improved road sign visibility.

*3) Surrounding complexity:* This parameter computes the amount of details that exist in the sign surroundings. It involves extracting the edges from all the surrounding areas and calculating the total number of edge pixels. The ratio between the number of edge pixels and the total number of pixels in these surrounding regions is employed to ascertain the shape complexity (C) of the road sign's surroundings, as follows:

$$C = \frac{N_E}{N_T} \qquad (6)$$

where, $N_E$ is the number of edge pixels in the surrounding regions and $N_T$ is the total number of pixels in these regions.

A complex environment around the road sign will result in a high complexity parameter value, leading to a reduced level of visibility. The overall complexity level of the surrounding regions of the sign will vary between high (0.2 - 1) and low (0 – 0.2) based on the value of the complexity parameter

*4) Occlusion:* This parameter quantifies the extent to which the road sign is partially obscured by objects like trees or leaves. It takes into account partial occlusion occurring on the top and right sides of the road sign region while disregarding occlusion on the left and bottom sides. The occlusion parameter (O) is formulated as follows:

$$O = 1 - \frac{A_O}{A_T} \qquad (7)$$

where, $A_O$ is the filled area of the apparent sign blob computed as the number of pixels and $A_T$ is the estimated area of road sign region computed as the bounding box area.

The level of occlusion can vary, being classified as either low (0 - 0.15) or high (0.15 - 1) based on the occlusion parameter value. Increased occlusion in the sign region would lead to reduced detectability and visibility of the road sign

*5) Tilting:* This parameter measures the tilting degree of road sign. The tilting parameter (T) is computed using the regionprops function on Python.

The degree of tilting can be categorized as either low (0 - 15°) or high (15° - 90°), contingent upon the tilting angle value. A pronounced tilt of the road sign would result in reduced detectability and consequently, a low visibility level.

*E. Determination of Visibility Levels*

Road signs are classified in this module using fuzzy logic in terms of visibility levels to: low, medium, or high. A Fuzzy Inference System (FIS) connecting parameters to the visibility level operates through a series of defined steps to determine the appropriate visibility label based on the input parameters. Considering the parameters calculated in the previous module (Readability, Color Difference, Surrounding Complexity, Occlusion, and Tilting) and their fuzzy sets mapped to visibility levels (Low, Medium, High), here's how the FIS functions:

*1) Input* variables and membership functions

- Parameters like Readability, Color Difference, Surrounding Complexity, Occlusion, and Tilting serve as input variables.

- Each parameter has fuzzy membership functions (e.g., low, high) that describe how input values correspond to these linguistic terms. These membership functions have defined ranges and shapes, such as triangular or Gaussian that assign degrees of membership to each linguistic term based on the input's value within its range. Table I shows the membership functions of the input parameters.

*2) Fuzzy* rules

- Based on expert knowledge or empirical data, fuzzy rules are established to connect the input parameters to the output visibility levels.

- For example, rules might state:

- "If Readability is High AND Occlusion is Low AND Color Difference is Low, THEN Visibility Level is High."

TABLE I. THE MEMBERSHIP FUNCTIONS OF THE FIVE VISIBILITY PARAMETERS

| Fuzzy Parameter | Membership Function Type | Parameter Range for Low | Membership Parameters for Low | Parameter Range for High | Membership Parameters for High |
|---|---|---|---|---|---|
| Readability | Triangular | 0 to 120 | a=0, b=60, c=120 | 120 to 255 | a=120, b=180, c=255 |
| Color difference | Triangular | 0 to 120 | a=0, b=60, c=120 | 120 to 255 | a=120, b=180, c=255 |
| Surrounding complexity | Triangular | 0 to 0.2 | a=0, b=0.1, c=0.2 | 0.2 to 1 | a=0.2, b=0.4, c=0.6 |
| Occlusion | Gaussian | 0 to .15 | Mean = 0.075 Standard Deviation = 0.0375 | 0.15 to 1 | Mean = 0.575 Standard Deviation = 0.2125 |
| Tilting | Trapezoidal | 0 to 15 | a=0, b=5, c=10, d=15 | 15 to 90 | a=15, b=20, c=85, d=90 |

*3) Inference* engine

- The inference engine evaluates the fuzzy rules based on the current input values.

- It calculates the degree to which each rule contributes to different visibility levels using fuzzy logic operations like AND, OR, and NOT.

- Aggregation methods, such as the Mamdani, combine the rules to determine the degree of support for each visibility level based on the input parameter values.

*4) Defuzzification*

- Once the inference engine processes the rules and combines their outputs, the defuzzification process aggregates the fuzzy output sets to derive a crisp, actionable output.

- This process converts the fuzzy output into a specific visibility level, such as Low, Medium, or High, based on methods like centroid, mean of maximum (MOM), or weighted average.

*5) Output* - determining visibility level

- The final step yields a specific visibility level determined by the FIS after processing the input parameters through the defined membership functions and rules.

- This output provides a clear and actionable visibility level based on the linguistic description or numerical range that best fits the input parameter combinations. Table II shows the membership functions of the output variable which is the visibility level.

The Fuzzy Inference System connects the input parameters related to visibility to the appropriate visibility level using fuzzy logic, allowing for a more understanding and decision-making process in scenarios where traditional binary or crisp logic might be insufficient.

The relationship between parameters and visibility levels is determined according to the following fuzzy rules:

*1) If* Readability is Low AND Occlusion is not high THEN Visibility Level is Low

*2) If* Occlusion is High THEN Visibility Level is Low

*3) If* Readability is High AND Occlusion is Low AND Color Difference is Low AND Surrounding Complexity is High AND Tilting is High THEN Visibility Level is Low

*4) If* Readability is High AND Occlusion is Low AND Color Difference is Low AND Surrounding Complexity is Low AND Tilting is High THEN Visibility Level is Medium

*5) If* Readability is High AND Occlusion is Low AND Color Difference is Low AND Surrounding Complexity is Low AND Tilting is Low THEN Visibility Level is High

*6) If* Readability is High AND Occlusion is Low AND Color Difference is Low AND Surrounding Complexity is Low AND Tilting is Low THEN Visibility Level is High

*7) If* Readability is High AND Occlusion is Low AND Color Difference is High AND Surrounding Complexity is Low AND Tilting is Low THEN Visibility Level is High

*8) If* Readability is High AND Occlusion is Low AND Color Difference is High AND Surrounding Complexity is High AND Tilting is Low THEN Visibility Level is High

*9) If* Readability is High AND Occlusion is Low AND Color Difference is High AND Surrounding Complexity is Low AND Tilting is High THEN Visibility Level is High

*10)If* Readability is High AND Occlusion is Low AND Color Difference is High AND Surrounding Complexity is High AND Tilting is High THEN Visibility Level is Medium

The rules connecting parameters to visibility levels in this fuzzy inference system have varying weights, indicating their significance in determining the visibility level. The lowest weight, at 0.2, is assigned to Rule 1, while Rule 2 holds a weight of 0.7, emphasizing the role of Occlusion in determining visibility. Rules 3 and 4, with weights of 0.8 and 0.9 respectively, highlight the combined impact of Readability, Occlusion, Color Difference, and Surrounding Complexity. Finally, Rules 5 to 10, each with a weight of 1.0, underscore the comprehensive consideration of Readability, Occlusion, Color Difference, Surrounding Complexity, and Tilting in determining the visibility level, demonstrating their paramount importance in decision-making.

TABLE II. THE MEMBERSHIP FUNCTIONS OF THE OUTPUT VARIABLE WHICH IS THE VISIBILITY LEVEL

| Visibility Level Fuzzy Output | Membership Function Type | Parameter Range for Low | Membership Parameters for Low | Parameter Range for Medium | Membership Parameters for Medium | Parameter Range for High | Membership Parameters for High |
|---|---|---|---|---|---|---|---|
| Low | Triangular | 0 to 0.33 | a=0, b=0.17, c=0.33 | 0.17 to 0.67 | a=0.17, b=0.42, c=0.67 | 0.33 to 1 | a=0.33, b=0.67, c=1 |
| Medium | Triangular | 0.17 to 0.67 | a=0.17, b=0.42, c=0.67 | 0.33 to 0.83 | a=0.33, b=0.58, c=0.83 | 0.67 to 0.83 | a=0.67, b=0.83, c=0.83 |
| High | Triangular | 0.33 to 1 | a=0.33, b=0.67, c=1 | 0.67 to 1 | a=0.67, b=0.83, c=1 | 0.67 to 1 | a=0.67, b=0.83, c=1 |

## IV. EXPERIMENTAL RESULTS

The visibility estimation system was tested on images of road signs taken by an in-vehicle camera in the United States. These in-vehicle images were obtained using a SAMSUNG ST65 camera, along with images from the VISAT™ Mobile Mapping System. All images were resized to 640x640 pixels. In this section, we will demonstrate the results of both the detection model and the visibility estimation model.

### A. Road Sign Detection Results

In this subsection, we evaluate the performance of the YOLOv8 detection model, which plays a crucial role in automatically identifying the road sign. The evaluation focuses on key metrics such as Accuracy, Precision, and mAP@0.5, providing insights into the model's accuracy and proficiency in object detection. The detection model underwent training for 150 epochs.

Fig. 5 presents a snapshot of quantitative metrics used to gauge the detection model's performance during training, including precision, recall, and mean average precision (mAP@0.5). These metrics shed light on the model's effectiveness in identifying road signs in in-vehicle images. Additionally, it is observed that box loss and class loss are converging.

The detailed breakdown of the model's performance is illustrated in the confusion matrix presented in Fig. 6, outlining true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN) for each object class (e.g., Warning sign, Regulatory sign, Stop sign, and background).

Precision measures the accuracy of true predictions made by the model. It is a crucial metric for object detection, as it assesses the model's ability to correctly identify objects without generating too many false positives. Recall assesses the model's ability to detect all relevant objects, reducing false negatives and ensuring no objects of interest are overlooked. mAP@0.5, a comprehensive metric, combines precision and recall, providing an aggregate evaluation of the model's performance across different object classes and considering precision-recall trade-offs.

The detection model has achieved remarkable results of the three performance metrics where mAP@0.5= 91.5%, Precision= 86.1%, and Recall= 90.5%. It is noticed that the model achieved better results of both Stop and Warning signs while missing some Regulatory signs. This happens because of the high effect of illumination on white signs especially when they are facing the sun.



Fig. 5. Performance metrics of the detection model throughout the training process for 150 epochs.



Fig. 6. Confusion matrix of the detection model.

### B. Visibility Estimation Results

Testing the fuzzy inference system (FIS) designed for road sign visibility demands a comprehensive procedure, especially when assessing its efficacy with a set of images. The initial step involves sourcing a diverse dataset of road sign images captured in various conditions, encompassing differing lighting, weather scenarios, angles, and distances. This dataset must cover a wide spectrum of potential real-world scenarios to ensure the FIS is tested under varying conditions. Upon gathering the images, preprocessing becomes pivotal. Standardizing the dataset involves resizing images to a uniform dimension, normalizing lighting conditions, and potentially applying filters or enhancements to accentuate visibility features present within the signs. Feature extraction follows, where specific visual features linked to the parameters considered in the FIS - such as readability, color difference, occlusion, surrounding complexity, and tilting - are identified and extracted from the images. This process is crucial to align the image data with the FIS parameters for subsequent analysis.

Integrating the FIS into the testing process involves applying the system to the extracted features from each image. This step aims to predict the visibility level for each sign based on the rules and weights defined within the system. Simultaneously, ground truth labeling becomes essential (assigning visibility labels to the images based on either human judgment or known visibility conditions captured during image acquisition). This establishes a benchmark against which the FIS's predictions can be evaluated. Post-prediction, an evaluation phase ensues where the predicted visibility levels are compared with the ground truth labels using various metrics, such as accuracy, precision and recall. Any discrepancies or misclassifications are carefully analyzed to understand potential shortcomings within the FIS. The procedure allows for iterative refinement. Any observed inconsistencies or errors guide adjustments to the FIS, such as tweaking membership functions, rules, or weights, aiming to enhance its accuracy and reliability.

*1) Evaluating* the Effectiveness of the Proposed Fuzzy Inference System: In evaluating the effectiveness of the proposed fuzzy inference system for visibility estimation, a distinct approach has been taken. Through a training phase involving 45 in-vehicle images, thresholds for detectability parameters were determined, crucial for classification into high, medium, or low visibility levels. This training set, comprising various road signs and diverse visibility scenarios, was instrumental in setting suitable threshold values, guided by expert decisions. These thresholds, rooted in the training phase, were then applied to a test set consisting of 50 in-vehicle images, including rectangular regulatory, diamond warning and stop signs, mirroring real-world diversity in visibility conditions.

The comparison between the decisions rendered by the proposed system and those of human experts unveils promising results. Out of 50 road signs tested, there was concurrence between the proposed system and expert judgments for 4 signs, representing an impressive 92% accuracy. Notably, even within the 4 instances of discordance, the disagreement usually amounted to merely one visibility level, showcasing a remarkable alignment between the proposed system's estimations and the human expert decisions. Fig. 7 shows examples of the proposed visibility estimation output along with the five visibility parameters and expert decision.

Further analysis revealed a nuanced performance difference in handling yellow and red versus white road signs. The system exhibited a higher proficiency with yellow and red signs owing to the impact of illumination on white color, affecting the accuracy of the color difference detectability parameter.

*2) Parametric influence on fuzzy inference system: shaping accuracy and decision dynamics:* The effectiveness and accuracy of the outcomes are profoundly influenced by the parameters incorporated within the system. These parameters, such as membership functions, threshold values, and rule weights, play a pivotal role in shaping the decisions

and predictions made by the system. Membership functions, serving as the backbone of fuzzy logic, define the degree of membership of an input to a specific linguistic term (like 'low,' or 'high'). Their design profoundly impacts the system's ability to interpret and categorize input data, significantly influencing the resulting output. Threshold values, especially in the context of detectability parameters for road sign visibility estimation, dictate the boundary between different visibility levels. Setting these thresholds involves a delicate balance; they need to be robust enough to delineate distinct visibility categories while remaining adaptable to varying environmental conditions.



Fig. 7. Visibility estimation outputs of the proposed system for road signs with expert decision: a) Low, b) High, c) Medium, d) Low.

Rule weights hold significance in the determination of the overall decision-making process within the fuzzy system. They assign importance or precedence to different rules, emphasizing the relative significance of specific parameters in contributing to the final output. Properly calibrated weights ensure that more critical parameters exert a more considerable influence on the system's decision. The effect of these parameters on the system's output is intricate and interconnected. Subtle adjustments or alterations in membership function shapes, threshold values, or rule weights can significantly impact the system's performance. Well-tuned parameters often lead to more accurate, reliable, and adaptable outcomes, enhancing the system's robustness across diverse datasets and real-world conditions.

Understanding the influence of these parameters allows for iterative refinement, facilitating continuous improvement in the system's accuracy and adaptability. Through careful

calibration and fine-tuning of these parameters, a fuzzy inference system can be optimized to yield more precise and dependable results, making it a valuable tool in addressing complex decision-making tasks where traditional binary logic falls short.

## V. CONCLUSIONS

In this work, we proposed a fully-automated system to detect road signs in the United States and estimate their visibility from images captured by an in-vehicle camera. Sign visibility is defined as the drivers' capability to perceive road signs on roadways, encompassing both the ability to detect the signs (Detectability) and the ability to read and recognize their contents (Readability).

The proposed system can be deployed in Driver Assistance Systems (DAS) or by transportation agencies. The proposed system has deployed YOLOv8 to build a detection model of three different road sign categories. Then, it measured five visibility parameters which are: readability of sign foreground, color difference between the sign and its surroundings, complexity of surroundings, sign occlusion, and sign tilting. The proposed system classifies road signs to three visibility levels: high, medium, and low. A Fuzzy Inference System (FIS) connecting these parameters to the visibility level operates through a series of defined steps to determine the appropriate visibility label based on the input parameters. The proposed system has achieved outstanding efficiency results with mAP@0.5= 91.5% for the detection model and an accuracy= 92% for the visibility estimation module. The accuracy of the proposed visibility estimation system has been compared with human expert pre-determined decisions.

The proposed system is distinguished by its being fully-automated, the efficiency of detecting road signs under various illumination and occlusion conditions, the ability to classify road signs visibility to multiple levels and the inclusion of both readability and detectability parameters of road signs from the perspective of driver.

In the future, we are planning to include more road sign categories in the visibility estimation system. Additionally, the size of dataset can be increased to improve the precision of the detection model. Hardware implementation can also be implemented based on the proposed computer vision system.

## REFERENCES

[1] H. B. Teklesenbet, N. H. Demoz, I. H. Jabiro, Y. R. Tesfay and E. Badidi, "Real-time Road Signs Detection and Recognition for Enhanced Road Safety," 2023 15th International Conference on Innovations in Information Technology, Al Ain, United Arab Emirates, pp. 132-137, 2023.

[2] K. Doman, D. Deguchi, T. Takahashi, Y. Mekada, I. Ide, H. Murase, Y. Tamatsu, "Estimation of traffic sign visibility toward smart driver assistance," in Intelligent Vehicles Symposium (IV), pp.45-50, June 2010.

[3] K. Doman, D. Deguchi, T. Takahashi, Y. Mekada, I. Ide, H. Murase, and U. Sakai, "Estimation of traffic sign visibility considering local and global features in a driving environment," In 2014 IEEE Intelligent Vehicles Symposium Proceedings, pp. 202-207, 2014.

[4] S. Maldonado-Bascón, S. Lafuente-Arroyo, P. Gil-Jiménez, H. Gómez-Moreno, and F. López-Ferreras, "Road-sign detection and recognition based on support vector machines," IEEE Trans. Intell. Transp. Syst., vol. 8, no. 2, pp. 264–278, Jun. 2007.

[5] A. de la Escalera, L. E. Moreno, M. A. Salichs, and J. M. Armingol, "Road traffic sign detection and classification," IEEE Transactions on Industrial Electronics, vol. 44, no. 6, pp. 848-859, Dec 1997.

[6] A. de la Escalera, J. MaArmingol, M. Mata, "Traffic sign recognition and analysis for intelligent vehicles," Image and Vision Computing, vol. 21, pp. 247–258, 2003.

[7] J. F. Khan, S. M. Bhuiyan, and R. R. Adhami, "Image segmentation and shape analysis for road-sign detection," IEEE Transactions on Intelligent Transportation Systems, vol.12, no.1, pp. 83-96, March 2011.

[8] X. Baro, S. Escalera, J. Vitria, O. Pujol, and P. Radeva, "Traffic sign recognition using evolutionary adaboost detection and forest-ECOC classification," IEEE Transactions on Intelligent Transportation Systems, vol. 10, pp. 113-126, 2009.

[9] J. Jiao, Z. Zheng, J. Park, Y. L. Murphey, and Y. Luo, "A robust multi-class traffic sign detection and classification system using asymmetric and symmetric features," IEEE International Conference on Systems, Man and Cybernetics, pp. 3421-3427, Oct. 2009.

[10] J. Abukhait, I. Abdel-Qader, J. S. Oh, and Abudayyeh, "Road sign detection and shape recognition invariant to sign defects," In 2012 IEEE International Conference on Electro/Information Technology (EIT), pp.1-6, May 2012.

[11] A. M. Mansour, J. Abukhait, and I. Zyout, "Speed sign recognition using independent component analysis," International Journal of Electrical, Electronics and Computer Systems (IJEECS), vol. 17, no. 01, pp. 832-838, Nov. 2013.

[12] J. Abukhait, I. Zyout, and A. M. Mansour, "Speed sign recognition using shape-based features," International Journal of Computer Applications (IJCA), vol. 84, no. 15, pp. 31-37, Dec. 2013.

[13] H. T. Ngoc, N. N. Vinh, N. T. Nguyen, and L. D. Quach, "Efficient Evaluation of SLAM Methods and Integration of Human Detection with YOLO Based on Multiple Optimization in ROS2," International Journal of Advanced Computer Science and Applications, vol. 14, no. 11, pp. 300-310, 2023.

[14] A. S. Sutikno, and R. Kusumaningrum, "Automated Detection of Driver and Passenger Without Seat Belt using YOLOv8," International Journal of Advanced Computer Science and Applications, vol. 14, no. 11, pp. 806-813, 2023.

[15] D. Tabernik and D. Skočaj, "Deep learning for large-scale traffic-sign detection and recognition," in IEEE Transactions on Intelligent Transportation Systems, vol. 21, no. 4, pp. 1427-1440, April 2020.

[16] A. A. Lima, M. M. Kabir, S. C. Das, M. N. Hasan, and M. F. Mridha, "Road sign detection using variants of yolo and r-cnn: An analysis from the perspective of Bangladesh," In Proceedings of the International Conference on Big Data, IoT, and Machine Learning, pp. 555-565, 2022..

[17] P. Kuppusamy, M. Sanjay, P. Deepashree, and C. Iwendi, "Traffic Sign Recognition for Autonomous Vehicle Using Optimized YOLOv7 and Convolutional Block Attention Module. Computers," Materials and Continua, vol. 77, no. 1, pp. 445-466, 2023.

[18] V. Balali, M. A. Sadeghi, and M. Golparvar-Fard, "Image-based retro-reflectivity measurement of traffic signs in day time," Elsevier Journal of Advanced Engineering Informatics, vol. 29, no. 4, pp. 1028-1040, 2015.

[19] F. Kimura, T. Takahashi, Y. Mekada, I. Ide, H. Murase, T. Miyahara, and Y. Tamatsu, "Measurement of visibility conditions toward smart driver assistance for traffic signals," in 2007 IEEE Intelligent Vehicles Symposium, pp. 636-641, June 2007.

[20] K. Doman, D. Deguchi, T. Takahashi, Y. Mekada, I. Ide, H. Murase, and Y. Tamatsu, "Estimation of traffic sign visibility considering temporal environmental changes for smart driver assistance," in 2011 IEEE Intelligent Vehicles Symposium (IV), pp. 667-672, June 2011.

[21] K. Mori, T. Kato, T. Takahashi, I. Ide, H. Murase, T. Miyahara, and Y. Tamatsu, "Visibility estimation in foggy conditions by in-vehicle camera and radar," Proc. International Conference on Innovative Computing, Information and Control, vol. 2, pp. 548-551, Aug. 2006.

[22] L. Simon, J. P. Tarel, and R. Br´emond, "Alerting the drivers about road signs with poor visual saliency," in Proc. 2009 IEEE Intelligent Vehicles Symposium, pp. 48–53, June 2009.

[23] J. Abukhait, I. Abdel-Qader, J. S. Oh, and O. Abudayyeh, "Occlusion-invariant tilt angle computation for automated road sign condition assessment," 2012 IEEE International Conference on Electro/Information Technology (EIT), pp. 1-6, 2012.

[24] J. Abukhait, "Visibility estimation of road signs considering detectability factors for driver assistance systems," WSEAS Transactions on Signal Processing, vol. 12, no. 1, pp. 111-117, 2014.

[25] S. Karami, and M. Taleai, "An innovative three-dimensional approach for visibility assessment of highway signs based on the simulation of traffic flow," Journal of Spatial Science, vol. 67, no. 2, pp. 203-218, 2022.

[26] J. Kulawik, M. Kubanek, and S. Garus, "The Verification of the Correct Visibility of Horizontal Road Signs Using Deep Learning and Computer Vision," Applied Sciences, vol. 13, no. 20, pp.11489, 2023.

[27] K. Doman, D. Deguchi, T. Takahashi, Y. Mekada, I. Ide, H. Murase, and U. Sakai, "Estimation of traffic sign visibility considering local and global features in a driving environment," In 2014 IEEE Intelligent Vehicles Symposium Proceedings, pp. 202-207, 2014.

[28] J. Abukhait, "Image-based Visibility Estimation of Road Signs in Cluttered Environment," International Journal of Computational and Applied Mathematics and Computer Science, vol. 2, pp.33-38, 2022.

[29] S. Zhang, C. Wang, M. Cheng, and J. Li, "Automated visibility field evaluation of traffic sign based on 3d Lidar point clouds," The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences, vol. 42, pp. 1185-1190, 2019.

[30] S. Zhang, C. Wang, L. Lin, C. Wen, C. Yang, Z. Zhang, and J. Li, "Automated visual recognizability evaluation of traffic sign based on 3D LiDAR point clouds," Remote Sensing, vol. 11, no. 12, pp.1453, 2019.

# Dual-Branch Grouping Multiscale Residual Embedding U-Net and Cross-Attention Fusion Networks for Hyperspectral Image Classification

Ning Ouyang, Chenyu Huang, Leping Lin

School of Information and Communication, Guilin University of Electronic Technology, Guilin, China

*Abstract*—Due to the high cost and time-consuming nature of acquiring labelled samples of hyperspectral data, classification of hyperspectral images with a small number of training samples has been an urgent problem. In recent years, U-Net can train the characteristics of high-precision models with a small amount of data, showing its good performance in small samples. To this end, this paper proposes a dual-branch grouping multiscale residual embedding U-Net and cross-attention fusion networks (DGMRU_CAF) for hyperspectral image classification is proposed. The network contains two branches, spatial GMRU and spectral GMRU, which can reduce the interference between the two types of features, spatial and spectral. In this case, each branch introduces U-Net and designs a grouped multiscale residual block (GMR), which can be used in spatial GMRUs to compensate for the loss of feature information caused by spatial features during down-sampling, and in spectral GMRUs to solve the problem of redundancy in spectral dimensions. Considering the effective fusion of spatial and spectral features between the two branches, the spatial-spectral cross-attention fusion (SSCAF) module is designed to enable the interactive fusion of spatial-spectral features. Experimental results on WHU-Hi-HanChuan and Pavia Center datasets shows the superiority of the method proposed in this paper.

*Keywords—U-Net; multiscale; cross-attention; hyperspectral image classification*

## I. INTRODUCTION

Hyperspectral images (HSI) include a wealth of spatial and spectral information [1], which can accurately characterize the physical attributes of features, enhance the ability to discriminate features, and bring great convenience to feature recognition. However, classification of hyperspectral images has its own special problems, such as the redundancy of information in spectral bands [2], the scarcity of training sample data [3], and class imbalance, which bring great challenges to hyperspectral image classification (HSIC).

Traditional HSIC classification methods, such as linear classifier [4], support vector machine [5] and random forest [6], can achieve good classification effect through improvement, but many original traditional methods rely on manual features, and the classification effect is poor when the number of samples is small and the HSI data dimension is high. Therefore, Principal Component Analysis (PCA) has been applied to HSIC by a large number of scholars [7-10]. By compressing the original data to reduce the spectral dimension, the information redundancy between bands and the possible

Hughes phenomenon can be avoided, which provides an effective treatment for subsequent feature extraction and enables the network to obtain higher classification accuracy.

With the development of deep learning, the encoder-decoder (U-net) [11] specially designed for biomedical image segmentation has been gradually applied to the field of hyperspectral image classification, which can obtain superior results with less training data. In the absence of datasets, Lin et al. [12] introduced U-Net to solve the problem of complex data capture in practice. Paul et al. [13] combined spectrum partitioning to reduce the redundancy of the spectrum, and then designed U-Net architecture by introducing deep separable convolution to reduce overfitting problems. Besides, due to the clear network structure of U-Net, any customized layer can be easily integrated into the existing network. For example, He et al. [14] embedded the Swin transformer into the classical CNN-based U-Net, which is dedicated to acquiring global contextual information of remote sensing images and obtaining deeper features in the master encoder. Xiao et al. [15] improved the spatial resolution of HSI by fusing spatial features of different scales and depths in the MSI for U-Net.

Moreover, in order to improve the classification performance of hyperspectral images, it has become a major research direction to jointly use spectral and spatial information to design classifiers. The construction of spatial and spectral information through dual branches can make full use of the information. Yang et al. [16] constructed a dual-channel CNN, extracting spectral and spatial information in each channel separately, and then connecting the spatial-spectral features by using cascade, but this simple feature connection cannot capture the complex relationship between the spatial-spectral features. Wang et al. [17] used the grouping strategy and the Long Short-Term Memory (LSTM) model to perceive spectral multi-scale information and obtain spatial context features in spectral finite element and spatial sub-network. Considering the different importance of spectral and spatial components, they used the method of adaptive feature combination for fusion. For effective fusion of spatial-spectral information, Sun et al. [18] designed a weighted self-attention fusion strategy, which combines the output weights of each branch of the previous network with the output weights of self-attention, and obtains efficient fusion on a multi-structured network. Yang et al. [19] used a dual-branch fusion mechanism to promote the exchange of feature information between the two branches through two upstream and downstream modules, so that local fine-grained features could be constructed in more detail and

global context information could be better utilized. These works provide new ideas for dual-branch feature extraction and fusion in HSIC.

In general, the method based on U-Net can better learn the representation of the input natural image, which is conducive to the classification of hyperspectral images to obtain high accuracy and obtain satisfactory results, but some small size information will be lost in the process of down-sampling. the design of fusion mechanism under two-branch conditions will also affect the effectiveness of the network. In this context, we propose a dual-branch grouping multiscale residual embedded U-Net and cross-attention fusion network. Among them, the main contributions are as follows:

- A dual-branching grouping multiscale residual embedded U-Net network (DGMRU) is proposed, which combines grouping multiscale residual block (GMR) and U-Net to extract rich global contextual information and deepen the function of feature network extraction.

- The grouping multiscale residual block (GMR) is constructed for digging multiscale information. The multiscale characteristics of this module enable the network to guide the network to focus on various types of samples at different scales, thereby improving the missed detection problem under spatial features and the redundancy problem under spectral features, and improving the effectiveness of feature extraction.

- A spatial-spectral cross-attention fusion module (SSCAF) is designed to cross-fuse the spatial and spectral features generated by the double branch, that is,

to fuse the parameters of the other branch into its own branch, increase the interaction of the two branches, and promote the full fusion of the two branches.

The rest of the paper is organized as follows. Section II describes the general framework of the DGMRU_CAF network, GMR and SSCAF, respectively. Section III discusses the dataset, the experimental settings, the experimental results and the discussion. Finally, in Section IV, conclusions are given.

## II. METHODOLOGY

### A. The Overall Framework of DGMRU_CAF

The DGMRU_CAF proposed in this paper is composed of DGMRU, SSCAF and classification network, as shown in Fig. 1. The DGMRU is divided into a spatial GMRU branch which takes the HSI neighbourhood block $p_n$ as input and a spectral GMRU branch which takes the spectral band $s_n$ as input. Each branch extracts corresponding features from the combined paths of U-Net and GMR with different nuclear scales, so as to obtain deeper feature information. In this regard, the designed GMR enhances the model's perception of multiscale spatial and spectral scales by grouping, multiscaling, and residual connection to retain more detailed feature information. Afterwards, in order to jointly utilize spatial and spectral information, the SSCAF module is constructed. Under the guidance of its own features, the module introduces the features of another branch and carries out interactive fusion to generate spatial-spectral features. Finally, in order to obtain the classification results of HSI, the obtained spatial-spectral features are passed through a classification network consisting of a fully connected layer and a softmax activation layer.



Fig. 1. The overall structure of DGMRU_CAF.

## B. The GMR Module

In this paper, a GMR module is proposed to retain more features without increasing parameters. For each branch, spatial GMR and spectral GMR are designed respectively.

*1) Spatial GMR:* As shown in Fig. 2(a), under the branch of spatial GMRU, for the intermediate features of the input space, spatial GMR uses the grouping module to group its spatial channels in sequence, so that each group of vectors contains different channel information, and each group of spatial vectors is expressed as:

$$c_{spa\_i} = [c_{i \times t}, c_{i \times t+1}, \cdots, c_{i \times t+t}], i = 0, \cdots, g-1 \quad (1)$$

where, $c_{i \times t}$ is the characteristic information corresponding to the channel in the $i \times t$ segment, t represents the number of channels in each group, and g represents the number of groups.

In the process of down-sampling, the features of small-size objects are easy to be weakened and lost, and it is difficult to recover these features by up-sampling, which leads to the misclassification of small-size objects. In order to solve this problem, this paper uses convolution of different sizes for multi-scale feature extraction after grouping to capture local features inherent in space. The convolution output of each group is:

$$s_i = W_i * x_i + b_i \quad (2)$$

where, $x_i$ is the feature vector of the ith group, $W_i$ is the weight coefficient of the ith group, and $b_i$ is the bias coefficient of the ith group.



(a)



(b)

Fig. 2. The structure of GMR module. (a) Spatial GMR (b) Spectral GMR.

Then, in order to complement the context information of features at different scales, all groups are merged by a cascade method. Finally, the rich low-frequency information is transmitted directly through the residual connection, which speeds up the training of the network. In short, using spatial GMR can extract more representative fine features.

*2) Spectral GMR:* Hyperspectral images contain a lot of spectral information, but the spectral information is redundant, which is easy to produce Hughes phenomenon and affect the classification results. In order to cope with this problem, and effectively capture the local relevant information of the spectral band. As shown in Fig. 2(b), for the spectral intermediate features, the grouping module of spectral GMR is used to group their spectral dimensions in sequence, so that each group of spectral vectors contains different spectral band information. Among them, the number of spectrum contained in each group and the distance between spectrum are related to the number of divided groups. Each set of spectral vectors is represented as:

$$r_{spe\_i} = [r_{i \times m}, r_{i \times m+1}, \cdots, r_{i \times m+m}], i = 0, \cdots, g-1 \quad (3)$$

Where, $r_{i \times m}$ is the characteristic information of the spectral dimension in the $i \times m$ segment, $m$ represents the number of spectral bands in each group, and g represents the number of groups.

After that, convolution of different scales is used to extract the grouped spectral features, so as to weaken the correlation between spectrums and reduce the redundancy of information. After that, the cascade method is used to merge the output spectral features of each group, which complements the local information of the spectral features of different scales and makes full use of the correlation between the spectral bands. Finally, the original global information is propagated directly by residual connection, which alleviates the problem of gradient degradation. In conclusion, the global and local information of spectra can be fully extracted by spectral GMR.

## C. The SSCAF Module

Considering the complementary characteristics between spatial and spectral features, in order to promote the effective fusion of these two types of features, a spatial-spectral cross-attention fusion (SSCAF) module is proposed in this paper. As shown in Fig. 3, the module is a combination of a cross self-attention module, a positional self-attention module (PAM) and a channel self-attention module (CAM). The cross self-attention operation is defined as follows:

$$y_{i,j} = \frac{1}{C(x)} \sum_{\forall i,j} f(f_i, f_j) g(f_i) \quad (4)$$

The $f_i, f_j$ represent the spatial and spectral feature vectors generated by the two branches, respectively, the function $f$ produces the adaptive weight vector between the two vectors, the function $g$ produces the feature representation of the input individual input vectors, and the normalization factor ss is defined as $C(X) = \sum_{\forall i,j} f(f_i, f_j)$.

To further establish internal connections, PAM and CAM modules are introduced to refine spatial and spectral features. Finally, the feature information is summed and

complementarily fused to obtain the final spatial-spectral fusion feature.



Fig. 3.   The structure of SSCAF module.

## III.   RESULTS

### A.  *Dataset and Experimental Setting*

In this section, to demonstrate the validity of the proposed method, we conduct a number of experiments on two datasets, which include WHU-Hi-HanChuan (HC) [20],[21] and Pavia Centre (PC). We divided the label samples in different ways for each data set. Table I and Table II provides the specific number of training, test sets and total samples for each class of the two data sets. The false color maps of the two datasets are shown in Fig. 4.



(b)

Fig. 4.   False color maps for the two datasets. (a) False color map of HC, (b) False color map of PC.

In the process of training the model, some parameters are set, where the training epoch is set to 200, the batch size is 16, the learning rate is 0.001, the weight decay is 1e-5, and the training is repeated 10 times for all the datasets. In order to prove the superiority of the proposed method, this paper conducts comparative experiments with six advanced methods, namely 2DCNN[22], SSRN[23], A2S2K[24], ASSMN[17], U-Net[11], HyperUnet[13]. The overall accuracy (OA), average accuracy (AA), Kappa coefficient and classification accuracy of single-class are used as the performance evaluation criteria of the model. The higher the each index, the better the classification effect will be.

TABLE I.          SAMPLE INFORMATION FOR EACH CLASS IN THE HC DATASET

| Class | Color | Class Name | Train | Test | Total |
|---|---|---|---|---|---|
| 1 | | Strawberry | 200 | 44535 | 44735 |
| 2 | | Cowpea | 200 | 22553 | 22753 |
| 3 | | Soybean | 200 | 10087 | 10287 |
| 4 | | Sorghum | 200 | 5153 | 5353 |
| 5 | | Water Spinach | 200 | 1000 | 1200 |
| 6 | | Watermelon | 200 | 4333 | 4533 |
| 7 | | Greens | 200 | 5703 | 5903 |
| 8 | | Trees | 200 | 17778 | 17978 |
| 9 | | Grass | 200 | 9269 | 9469 |
| 10 | | Red Roof | 200 | 10316 | 10516 |
| 11 | | Gray Roof | 200 | 16711 | 16911 |
| 12 | | Plastic | 200 | 3479 | 3679 |
| 13 | | Bare Soil | 200 | 8916 | 9116 |
| 14 | | Road | 200 | 18360 | 18560 |
| 15 | | Bright Object | 200 | 936 | 1136 |
| 16 | | Water | 200 | 75201 | 75401 |
| Total | | | 3200 | 254330 | 257530 |

TABLE II.          SAMPLE INFORMATION FOR EACH CLASS IN THE PC DATASET

| Class | Color | Class Name | Train | Test | Total |
|---|---|---|---|---|---|
| 1 | | Water | 82 | 742 | 824 |
| 2 | | Trees | 82 | 738 | 820 |
| 3 | | Asphalt | 81 | 735 | 816 |
| 4 | | Self-Blocking Bricks | 80 | 728 | 808 |
| 5 | | Bitumen | 80 | 728 | 808 |
| 6 | | Tiles | 100 | 1160 | 1260 |
| 7 | | Shadows | 47 | 429 | 476 |
| 8 | | Meadows | 82 | 742 | 824 |
| 9 | | Bare Soil | 82 | 738 | 820 |
| Total | | | 716 | 6740 | 7456 |

### B.  *Analysis of Classification Results of Dataset*

- Classification Maps and Result of HC Dataset.

The results on the HC dataset are shown in Table III, with the best OA, AA, and Kappa results highlighted in bold. Fig. 5 shows the classification diagram for the different methods.

It can be seen from Table III that our method achieves the best performance, with OA of 96.22%, AA of 96.62%, and Kappa of 95.57%. Compared with other methods, OA, AA, and Kappa are increased by at least 0.8%, 0.76%, and 0.92%. This is because the proposed method has the new idea of combining dual-branch and U-Net, which improves the ability of convolutional feature extraction, so that the method in this paper can achieve the best performance. The grouping multiscale residual block is designed to extract features with different kernel sizes in each group, and reduce the loss of feature information to construct effective feature extraction. The classification results of HSI prove the validity of the method. In addition, it can be seen that the OA of 2DCNN is

the lowest, only 78.91%, which is because 2DCNN is trained only on the spatial dimension, ignoring the information between spectrum, and the model performance is poor. Compared with 2DCNN, U-Net constructs U-shaped network structure, improves classification accuracy and performance, and improves 13.38%, 15.84% and 15.3% respectively in the three evaluation criteria. HyperUnet networks, which combine U-Net and grouping ideas, perform poorly on this dataset, possibly because of poor adaptability to large datasets. In addition, it can be observed that the evaluation value of SSRN is comparable to that of U-Net. SSRN extracts spatial-spectral features through the combination of two continuous spectral blocks and spatial blocks. However, the input of the spatial block comes from the spectral block, which leads to the loss of some spatial information in the spectral block, resulting in poor classification accuracy. The OA of A2S2K is better than that of SSRN, increased by 1.69%, which indicates that the introduction of attention mechanism and adaptive methods has a significant impact on the network. Compared with other single-branch algorithms, the dual-branch ASSMN results in better OA values, which indicates that the full use of spatial

and spectral feature information can achieve superior classification results, and the effect is much better than that of single spatial or spectral information. Although the effectiveness of the method in this paper is inferior to other algorithms in some categories, the results of these methods are very close to the results of the best classification, so the OA, AA and kappa coefficients of the method in this paper are the highest among these methods.

From the classification diagram shown in Fig. 5, the "salt and pepper" noise is the most severe because spectral information is not included in 2DCNN, while the classification diagram of other networks shows stronger classification ability because spectral information is taken into account. The method proposed in this paper considers the spatial features of different scales and solves the redundancy problem to obtain more small-size objects and feature information. Therefore, for classification maps with more small sizes, the method proposed in this paper is easier to obtain more accurate and cleaner classification maps, and the classification results of various categories correspond to the results in Table III.



|  (a)  |  (b)  |  (c)  |  (d)  |  (e)  |  (f)  |  (g)  |  (h)  |

Fig. 5. Classification maps of different methods in HC dataset, (a) Ground truth, (b) 2DCNN, (c) SSRN, (d) A2S2K, (e) U-Net, (f) HyperUnet, (g) ASSMN, (h) Proposed.

TABLE III. CLASSIFICATION RESULTS OF THE HC DATASET

| Class | 2DCNN | SSRN | A2S2K | U-Net | HyperUnet | ASSMN | Proposed |
|---|---|---|---|---|---|---|---|
| 1 | 90.78 | 92.71 | 92.54 | 93.27 | 73.72 | **95.80** | 95.12 |
| 2 | 84.81 | 91.65 | 86.2 | 94.52 | 73.59 | **95.80** | 93.38 |
| 3 | 91.43 | 95.50 | **97.31** | 94.65 | 74.38 | 95.07 | 97.04 |
| 4 | 97.49 | **99.20** | 99.45 | 99.45 | 92.85 | 99.10 | 98.91 |
| 5 | 98.70 | **100** | **100** | **100** | 92.70 | **100** | **100** |
| 6 | 70.69 | 95.5 | 90.53 | 95.19 | 66.74 | **97.71** | 95.47 |
| 7 | 49.04 | 96.35 | 96.82 | **99.35** | 84.69 | 95.75 | 99.10 |
| 8 | 63.24 | 89.83 | 82.07 | 9.17 | 61.84 | 88.77 | **96.38** |
| 9 | 63.02 | 92.08 | **95.46** | 91.74 | 64.89 | 97.85 | 94.47 |
| 10 | 99.46 | 97.37 | **98.92** | 90.25 | 90.92 | 92.54 | 95.86 |
| 11 | 47.87 | 90.88 | **99.15** | 86.80 | 83.32 | 94.57 | 89.52 |
| 12 | 75.88 | 98.93 | 97.64 | 98.47 | 65.24 | **100** | 99.97 |
| 13 | 59.52 | 78.88 | 91.07 | 87.93 | 70.59 | 87.78 | **95.33** |
| 14 | 82.55 | 94.03 | 91.67 | 93.66 | 76.28 | **97.18** | 96.96 |
| 15 | 97.54 | **100** | 99.57 | 97.75 | 84.93 | 99.14 | 99.78 |
| 16 | 80.70 | 92.81 | **98.73** | 90.85 | 94.85 | 96.74 | 98.61 |
| OA（%） | 78.91 | 92.61 | 94.30 | 92.29 | 80.75 | 95.42 | **96.22** |
| AA（%） | 78.29 | 94.11 | 94.83 | 94.13 | 78.22 | 95.86 | **96.62** |
| Kappa×100 | 75.70 | 91.40 | 93.33 | 91.00 | 77.69 | 94.65 | **95.57** |

- Classification Maps and Result of PC Dataset

The results on the PC dataset are shown in Table IV, with the best OA, AA, and Kappa results highlighted in bold. Fig. 6 shows the classification diagram for the different methods.

As shown in Table IV, on this PC dataset, all methods, including 2DCNN, achieve decent classification results. Obviously, the AA of both 2DCNN and SSRN are lower than 90%, which is due to poor classification accuracy in some categories, and the classification accuracy of some categories is less than 80%. The values of A2S3K, U-Net and HyperUnet in OA, AA and Kappa all reach more than 90%, but it is still difficult to improve the classification accuracy for some categories. As a two-branch multi-scale network, ASSMN has more stable classification results. Method of this paper is superior, has the best OA, AA and Kappa evaluation values, and achieves the best accuracy for some specific categories, such as Class 4 Self-Blocking Bricks and Class 5 Bitumen, which further proves its validity in terrain classification.

As shown in Fig. 6, our method is smoother and more consistent.

## C. Ablation Analysis

In this part, extensive ablation experiments are conducted to demonstrate the validity of the proposed GMR, SSCAF on the two datasets.

The validity analysis of GMR is shown in Table V. It can be seen that without GMR, the values of OA, AA and Kappa of the model are the lowest in the experiment, because it will lead to some small-size samples being ignored in the process of down-sampling. In contrast, the simultaneous presence of GMR modules with two branches can extract spectral and spatial features more effectively, which contributes to the final classification, and its OA, AA, and Kappa can achieve the best results compared with other comparison strategies. Among them, OA increased by 2.6% and 1.51% in the two datasets, respectively, which means the necessity of GMR. In addition, the OA value of "Only Spe-GMR" is higher than that of "Only Spa-GMR", because the HSI contains enough spectral information to extract more useful feature information from it.

The results of SSCAF ablation experiments are shown in Table VI. It can be found that the integration of SSCAF into the two branches of "With GMR" has significantly improved network performance, which means that SSCAF can complement each other with spatial and spectral information to contribute to the final classification decision. Compared with without SSCAF, OA is increased by 4.54% and 3.03%, AA is increased by 3.14% and 3.59%, and Kappa is increased by 5.32 and 3.46, respectively, which fully proves the necessity of the existence of SSCAF.

TABLE IV.    CLASSIFICATION RESULTS OF THE PC DATASET

| Class | 2DCNN | SSRN | A2S2K | U-Net | HyperUnet | ASSMN | Proposed |
|---|---|---|---|---|---|---|---|
| 1 | **100** | 99.77 | **100** | 98.78 | 99.88 | 99.99 | 99.82 |
| 2 | 91.41 | 97.54 | 64.51 | 74.29 | 91.77 | **95.95** | 90.42 |
| 3 | 96.48 | 90.63 | **99.89** | 86.12 | 93.07 | 92.27 | 90.3 |
| 4 | 95.04 | 99.57 | 97.71 | 99.76 | 96.47 | **100** | **100** |
| 5 | 88.54 | 40.29 | 97.56 | 96.94 | 94.30 | 92.82 | **99.09** |
| 6 | 67.07 | 99.14 | **99.95** | 82.65 | 85.05 | 96.55 | 95.74 |
| 7 | 90.05 | 77.97 | 95.71 | **98.35** | 88.24 | 91.45 | 97.16 |
| 8 | 98.75 | 98.14 | 98.34 | 95.68 | **99.65** | 96.50 | 97.78 |
| 9 | 7.49 | **100** | **100** | 87.22 | 99.20 | 95.51 | 97.86 |
| OA（%） | 94.27 | 95.28 | 97.35 | 95.07 | 97.45 | 97.64 | **98.11** |
| AA（%） | 81.65 | 89.23 | 94.85 | 91.09 | 94.18 | 96.01 | **96.48** |
| Kappa×100 | 91.73 | 93.32 | 96.24 | 93.04 | 96.38 | 96.66 | **97.32** |



| (a) | (b) | (c) | (d) | (e) | (f) | (g) | (h) |

Fig. 6.    Classification maps of different methods in PC dataset,(a) Ground truth, (b) 2DCNN, (c) SSRN, (d)A2S2K, (e) U-Net, (f) HyperUnet, (g) ASSMN, (h) Proposed.

TABLE V.    EFFECTIVENESS ANALYSIS OF GMR

| Strategy | HC | | | PC | | |
|---|---|---|---|---|---|---|
| | OA | AA | Kappa | OA | AA | Kappa |
| Without GMR | 93.62 | 95.39 | 92.58 | 96.60 | 92.03 | 93.91 |
| Only Spe-GMR | 95.39 | 95.71 | 94.33 | 96.63 | 94.33 | 95.18 |
| Only Spa-GMR | 94.44 | 95.43 | 93.64 | 96.36 | 93.63 | 94.86 |
| With GMR | **96.22** | **96.62** | **95.57** | **98.11** | **96.48** | **97.32** |

TABLE VI. EFFECTIVENESS ANALYSIS OF SSCAF

| Strategy | HC | | | PC | | |
|---|---|---|---|---|---|---|
| | OA | AA | Kappa | OA | AA | Kappa |
| Without SSCAF | 91.68 | 93.48 | 90.25 | 95.08 | 92.89 | 93.86 |
| With SSCAF | **96.22** | **96.62** | **95.57** | **98.11** | **96.48** | **97.32** |

*D. Discussion of Training Times and Testing Times*

In order to measure the efficiency of the proposed method, this paper conducted comparative experiments in training and testing time, and the results are shown in Table VII. 2DCNN has the least training time and testing time than other methods, because the simple 2DCNN architecture has fewer training parameters, but the classification accuracy is relatively low. Both U-net and HyperUnet have encoding and decoding path modules, and the addition of more convolutional layers makes the consumption time slightly longer than that of 2DCNN. SSRN and A2S2K use 3D convolution and introduce ResNet, which speeds up convergence and reduces training time. In

ASSMN, the combination of dual-branch and multi-scale, together with its strategy of spectrum grouping and spatial grouping, makes the model more complex and requires longer training and testing time. However, the training time and testing time of the method in this paper are average among these comparison methods. The attention mechanism used in the SSCAF module increases the complexity of the proposed network, and the obtained training time and testing time are not the shortest. However, the method proposed in this paper can strike a good balance between accuracy and efficiency, and has certain advantages.

TABLE VII. RUNNING TIME OF DIFFERENT METHODS ON TWO DATASETS

| Dataset | | 2DCNN | SSRN | A2S2K | U-Net | HyperUnet | ASSMN | Proposed |
|---|---|---|---|---|---|---|---|---|
| HC | Train(s) | 8.96 | 48.53 | 63.96 | 22.54 | 38.05 | 192.15 | 82.46 |
| | Test(s) | 42.15 | 123.80 | 205.50 | 61.27 | 82.04 | 121.01 | 138.05 |
| PC | Train(s) | 3.11 | 7.84 | 11.75 | 7.45 | 8.65 | 54.90 | 18.88 |
| | Test(s) | 9.32 | 18.74 | 23.53 | 34.74 | 23.30 | 69.78 | 73.41 |

## IV. CONCLUSION

In this paper, we propose a dual-branch grouping multiscale residual embedded U-Net and cross-attention fusion network for hyperspectral image classification to improve the classification accuracy in the presence of sparse training samples. The designed DGMRU module is used to extract multiscale context information feature, which is suitable for the case of insufficient HSI samples. Among them, the designed GMR module increases the receptive field without adding parameters, and the feature extraction effect is better than that of the non-existent GMR module, which proves the necessity of this module. In addition, the proposed SSCAF maximizes the utilization of spatial-spectral features by constructing the intrinsic relationship between spatial and spectral features through cross-attention. Compared with other advanced algorithms, the method proposed in this paper has the best experimental results, and in the two data sets, OA increases by 0.8% and 0.47% at least, which is feasible and effective. In the future, we will consider further reducing the complexity of the network model and improving the computational efficiency while maintaining the classification accuracy.

## REFERENCES

[1] X. Li, Z. Li, H. Qiu, G. Hou, and P. Fan, "An overview of hyperspectral image feature extraction, classification methods and the methods based on small samples," Appl. Spectrosc. Rev., vol. 58, no. 6, pp. 367–400, Jul. 2023.

[2] R. N. Patro, S. Subudhi, P. K. Biswal, and F. Dell'acqua, "A Review of Unsupervised Band Selection Techniques: Land Cover Classification for Hyperspectral Earth Observation Data," IEEE Geosci. Remote Sens. Mag., vol. 9, no. 3, pp. 72–111, Sep. 2021.

[3] X. Wang, J. Liu, W. Chi, W. Wang, and Y. Ni, "Advances in Hyperspectral Image Classification Methods with Small Samples: A Review," Remote Sens., vol. 15, no. 15, Art. no. 15, Jan. 2023.

[4] M. Shambulinga and G. Sadashivappa, "Supervised hyperspectral image classification using svm and linear discriminant analysis," Int. J. Adv. Comput. Sci. Appl., vol. 11, no. 10, pp. 403–409, 2020.

[5] J. Dr. J. H. Harikiran, "Hyperspectral image classification using support vector machines," IAES Int. J. Artif. Intell. IJ-AI, vol. 9, no. 4, p. 684, Dec. 2020.

[6] F. Tong and Y. Zhang, "Exploiting Spectral–Spatial Information Using Deep Random Forest for Hyperspectral Imagery Classification," IEEE Geosci. Remote Sens. Lett., vol. 19, pp. 1–5, 2022.

[7] M. P. Uddin, M. A. Mamun, and M. A. Hossain, "PCA-based feature reduction for hyperspectral remote sensing image classification," IETE Tech. Rev., vol. 38, no. 4, pp. 377–396, 2021.

[8] H. Fu, G. Sun, J. Ren, A. Zhang, and X. Jia, "Fusion of PCA and Segmented-PCA Domain Multiscale 2-D-SSA for Effective Spectral-Spatial Feature Extraction and Data Classification in Hyperspectral Imagery," IEEE Trans. Geosci. Remote Sens., vol. 60, pp. 1–14, 2022.

[9] Q. Liu, D. Xue, Y. Tang, Y. Zhao, J. Ren, and H. Sun, "PSSA: PCA-Domain Superpixelwise Singular Spectral Analysis for Unsupervised Hyperspectral Image Classification," Remote Sens., vol. 15, no. 4, p. 890, 2023.

[10] X. Zhang, X. Jiang, J. Jiang, Y. Zhang, X. Liu, and Z. Cai, "Spectral–Spatial and Superpixelwise PCA for Unsupervised Feature Extraction of Hyperspectral Imagery," IEEE Trans. Geosci. Remote Sens., vol. 60, pp. 1–10, 2022.

[11] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, Springer, 2015, pp. 234–241.

[12] M. Lin, W. Jing, D. Di, G. Chen, and H. Song, "Context-aware attentional graph U-Net for hyperspectral image classification," IEEE Geosci. Remote Sens. Lett., vol. 19, pp. 1–5, 2021.

[13] A. Paul and S. Bhoumik, "Classification of hyperspectral imagery using spectrally partitioned HyperUnet," Neural Comput. Appl., pp. 1–10, 2022.

[14] X. He, Y. Zhou, J. Zhao, D. Zhang, R. Yao, and Y. Xue, "Swin transformer embedding UNet for remote sensing image semantic segmentation," IEEE Trans. Geosci. Remote Sens., vol. 60, pp. 1–15, 2022.

[15] J. Xiao, J. Li, Q. Yuan, and L. Zhang, "A dual-UNet with multistage details injection for hyperspectral image fusion," IEEE Trans. Geosci. Remote Sens., vol. 60, pp. 1–13, 2021.

[16] J. Yang, Y. Zhao, J. C.-W. Chan, and C. Yi, "Hyperspectral image classification using two-channel deep convolutional neural network," in 2016 IEEE international geoscience and remote sensing symposium (IGARSS), IEEE, 2016, pp. 5079–5082.

[17] D. Wang, B. Du, L. Zhang, and Y. Xu, "Adaptive spectral–spatial multiscale contextual feature extraction for hyperspectral image classification," IEEE Trans. Geosci. Remote Sens., vol. 59, no. 3, pp. 2461–2477, 2020.

[18] L. Sun, Y. Fang, Y. Chen, W. Huang, Z. Wu, and B. Jeon, "Multi-structure KELM with attention fusion strategy for hyperspectral image classification," IEEE Trans. Geosci. Remote Sens., vol. 60, pp. 1–17, 2022.

[19] L. Yang et al., "FusionNet: a convolution–transformer fusion network for hyperspectral image classification," Remote Sens., vol. 14, no. 16, p. 4066, 2022.

[20] Y. Zhong, X. Hu, C. Luo, X. Wang, J. Zhao, and L. Zhang, "WHU-Hi: UAV-borne hyperspectral with high spatial resolution (H2) benchmark datasets and classifier for precise crop identification based on deep convolutional neural network with CRF," Remote Sens. Environ., vol. 250, p. 112012, 2020.

[21] Y. Zhong et al., "Mini-UAV-borne hyperspectral remote sensing: From observation and processing to applications," IEEE Geosci. Remote Sens. Mag., vol. 6, no. 4, pp. 46–62, 2018.

[22] X. Yang, Y. Ye, X. Li, R. Y. Lau, X. Zhang, and X. Huang, "Hyperspectral image classification with deep learning models," IEEE Trans. Geosci. Remote Sens., vol. 56, no. 9, pp. 5408–5423, 2018.

[23] Z. Zhong, J. Li, Z. Luo, and M. Chapman, "Spectral–spatial residual network for hyperspectral image classification: A 3-D deep learning framework," IEEE Trans. Geosci. Remote Sens., vol. 56, no. 2, pp. 847–858, 2017.

[24] S. K. Roy, S. Manna, T. Song, and L. Bruzzone, "Attention-based adaptive spectral–spatial kernel ResNet for hyperspectral image classification," IEEE Trans. Geosci. Remote Sens., vol. 59, no. 9, pp. 7831–7843, 2020.

# FPGA-based Implementation of a Resource-Efficient UNET Model for Brain Tumour Segmentation

Modise Kagiso Neiso[1], Dr. Nicasio Maguu Muchuka[2], Dr. Shadrack Maina Mambo[3]

Department of Electrical & Electronics Engineering, PAUSTI, Juja, Kenya[1, 2, 3]
Department of Electrical & Control Engineering, Egerton University, Nakuru, Kenya[2]
Electrical Engineering Department, Walter Sisulu University, Ibika, South Africa[3]

*Abstract*—In this study an optimized UNET model is used for FPGA-based inference in the context of brain tumour segmentation using the BraTS dataset. The presented model features reduced depth and fewer filters, tailored to enhance efficiency on FPGA hardware. The implementation leverages High-Level Synthesis for Machine Learning (HLS4ML) to optimize and convert a Keras-based UNET model to Hardware Description Language (HDL) in the Kintex Ultrascale (xcku085-flva1517-3-e) FPGA. Resource strategy, First in First out (FIFO) depth optimization, and precision adjustment were employed to optimize FPGA resource utilization. Resource strategy is demonstrated to be effective, with resource utilization reaching a saturation point at a 1000-reuse factor. Following FIFO optimization, significant reductions are observed, including a 55 percent decrease in Block RAM (BRAM) usage, a 43 percent reduction in Flip-Flops (FF), and a 49 percent reduction in Look-Up Tables (LUT). In C/RTL co-simulation, the proposed FPGA-based UNET model achieves an Intersection over Union (IoU) score of 74 percent, demonstrating comparable segmentation accuracy to the original Keras model. These findings underscore the viability of the optimized UNET model for efficient brain tumour segmentation on FPGA platforms.

*Keywords—UNET; field programmable gate array; high-level synthesis for machine learning; brain tumour segmentation*

## I. INTRODUCTION

Brain tumours are abnormal growths of cells in or around the brain and can be cancerous or non-cancerous [1]. Mutations in the DNA, radiation exposure and immune system problems are the causes of brain tumours which cause a noticeable mortality and low recovery rates. Early detection of brain tumours is crucial as it increases the possibility of successful treatment [2]. Magnetic resonance imaging (MRI) and computed tomography (CT) are imaging modalities used to diagnose brain tumours, with MRI producing more detailed brain scans [3]. Diagnosis of brain tumours from MRI requires skilled manpower in the medical field [4]. The expertise required for brain tumour diagnosis is insufficient and is susceptible to the human error factor, which has resulted in the implementation of deep learning (DL) to predict tumours to assist doctors [5].

Convolutional Neural Networks (CNN) are the deep DL model that perform better for feature extraction, but they require large datasets for efficient training which is hindrance for applications in medical imaging as large dataset are not easily accessible. Ronnerberger et al. proposes a UNET architecture that requires less image samples for successful model learning [6]. The UNET architecture has a drawback of consuming a lot of resources and computing inefficiency when applied in CPU and GPU [7]. The computational inefficiency of existing tumour prediction methods has become a critical concern in the medical imaging field. In response to this challenge, recent research has focused on areas to apply the UNET model for brain tumour detection in field programmable gate arrays (FPGA's), due to their inherent speed advantage over traditional processors. In a related study [8], the FPGA implementation of CNN was discussed, emphasizing on the necessity of a balanced design that considers resource utilization against performance.

In this work the challenge of balancing computational resource against performance in UNET for brain tumour prediction is addressed by proposing a modified UNET model and a comprehensive optimization of hardware resource utilization during FPGA inference. The modifications methods used in this work are tailored to enhance efficiency and efficacy of the UNET model for brain tumour prediction. The UNET model was optimized by reducing the size of original UNET model and FPGA hardware optimization strategies used entailed resource strategy, FIFO buffer depth optimization and precision.

The rest of this paper is organized as follows: Section II is literature review of relevant theory and related work; Section III describes materials and methods to carry out the work; Section IV is analysis and interpretation of results and discussion and conclusion is given in Section V and Section VI respectively.

## II. LITERATURE REVIEW

Brain tumour segmentation is the technique of automatically detecting and labelling malignant brain tissues depending on tumour type [9]. Convolutional Neural Network has been successful in image segmentation applications in deep learning with applications in the medical field. The use of CNN in image segmentation entails deciding on the dataset to train the model; CNN architecture to use; loss function and the back-propagation weight adjustment algorithm.

### A. Brain Tumour Dataset

Brain MRI scans is private patient's data which require confidential handling by health practitioners. However with the rise in the use of technology in the medical field, there is a need for the data to be made public and anonymous of patient's identity [10]. International research institutes have

made this data available for public use in developing medical imaging solutions such as The Cancer Genome Atlas (TCGA) dataset, BrainWeb dataset and MICCAI Brain Tumour Segmentation (BraTS 2020) Challenge dataset.

*1) The Cancer Genome Atlas (TCGA):* TCGA has different data such as TCGA Lower Grade Gliomas (LGG) and Glioblastoma Multiforme (GBM) which are basically grades of gliomas standardized by the World Health Organization. LGG are less aggressive while GBM is an exceptionally aggressive kind of brain glioma that develops from astrocytes or their progenitors [11]. In reserach [12], automatic and manual identification of GBM sub-compartments segmentation was performed and results showed that automated segmentation gave the highest area under the curve (AUC) as compared to manual segmentation. The TCGA dataset was created to identify a causal relation between genomic alterations and cancers [13]. TCGA data does not have specific pixel segmentation of tumour regions which will require additional radiologist expert knowledge to label the scans for training a machine learning model [14].

*2) BrainWeb:* Simulated Brain Database (SBD): The SBD currently has brain MRI data that has been simulated using two anatomical models: normal and multiple sclerosis (MS). T1-weighted (spin-lattice relaxation), T2-weighted (spin-spin relaxation), and PD-weighted (proton-density) sequences were employed to simulate entire three-dimensional data volumes for these models [15]. In study [16], the SBD database was utilized for improving the magnetic resonance imaging (MRI) segmentation using fuzzy C-means method and obtained experimental results that were more stable and accurate when compared to existing methods. SBD was created for computer aided image analysis by providing samples with ground truth. The SBD dataset is simulated for general use in computer-based analysis algorithms which diminishes its appeal when compared to other datasets sourced from actual patients.

*3) MICCAI Brain Tumour Segmentation (BraTS 2020) challenge dataset:* The BraTS 2020 dataset is made up of clinically obtained pre-operative multimodal MRI images of Glioblastoma (GBM/HGG) and lower grade glioma (LGG) that were collected from multiple institutions [17,18]. It contains a diagnosis that has been verified pathologically and has been divided into training and validation data. The dataset has expert manual segmentations that define the boundaries of the tumour regions, which include enhancing tumour, the peritumoral edema, and the necrotic and non-enhancing tumour core. All the scans of BraTS are in Neuroimaging Informatics Technology Initiative (NIfTIS) file format and they describe native (T1) and post-contrast T1-weighted (T1Gd), T2-weighted (T2), and T2-Fluid Attenuated Inversion Recovery (FLAIR) volumes, and were acquired using various clinical protocols and scanners from numerous institutions. The BraTS dataset has been used in several research with [19] introducing a residual mobile U-Net (RMU-Net) for MRI brain tumour segmentation. The RMU-Net archived dice coefficient scores for WT, TC, and ET on the BraTS 2020 dataset of 91.35%, 88.13%, and 83.26%, respectively, and 91.76%, 91.23%, and 83.19% on the BraTS 2019 dataset, and 90.80%, 86.75%, and 79.36% on the BraTS 2018 dataset.

The BraTS dataset is more versatile than the TCGA as it has manual segmentation that has been done by experts, which is essential for the training of a prediction model. In comparison to the SBD, BraTS dataset is more competitive as it has been sourced from real MRI scans that represents the real life cases as compared to simulated SBD. BraTS dataset is also specialized for segmentation with UNET as it has ground truth that has multiclass labels [20].

*B. U-NET*

U-NET is a convolutional neural network that was proposed by [6] for biomedical image segmentation applications. It was optimized to archive accurate prediction with few training images, as they are normally few sample images in the medical field. The original model constitutes of a contraction path which records context and expansion path which enables accurate localization. The UNET model is computationally demanding which translates into high resource utilization on hardware. There have been proposed methods to reduce the high resource demand by the UNET.

In study [21] a reduced UNET model architecture for classification of weeds and crops using segmentation was proposed. Reduction of the model was archived by reducing the number of filters per convolution layer. The proposed model in [21] has parameters that are 27% smaller while maintaining accuracy of 95% and an error rate that is 7% lower than the original UNET model. The reduction is however done on the number of filter per layer, and maintains the UNET architecture in terms of layers. In [22], the reduced depth UNET architecture with three down-sampling and two up-sampling sections is proposed, replacing the five down-sampling and four up-sampling sections of the original UNET architecture. Results obtained showed that the approach produced more accurate results [22]. There is gap in the existing work to combine reduction in number of filters and model depth as proposed by [21, 22].

*1) Evaluation criteria:* In order to evaluate segmentation there are metrics in image processing that can be used for quality assurance. The output image pixels are compared against those of the ground truth to establish the extent of accuracy. Intersection over union (IoU) and dice similarity coefficient (DSC) are the most common used metrics in medical imaging.

*a) Intersection over union (Iou):* Intersection over union is an evaluation metric that measures the intersection between the predicted mask and the actual mask, also known as the Jaccard index [23]. IoU calculation incorporates two indicators of false positive (FP) and true positive (TP) results. A true positive occurs when the model accurately predicts that a pixel is a component of an object when in fact it is. If the model forecasts a pixel as belonging to an item when in fact it belongs to the background, this is known as a false positive. The intersection of the anticipated segmentation mask and the ground truth mask, when compared to the union of the two

masks, is known as the IoU. IoU is particularly useful in multiclass segmentation where there is an imbalance in classes.

$$IoU = TP / (TP + FP + FN) \qquad (1)$$

where, TP is the number of true positives, FP is the number of false positives, and FN is the number of false negatives.

*b) Dice similarity coefficient:* The Dice similarity coefficient also termed as Srensen-Dice index or simply the Dice coefficient is a statistical instrument used to determine the similarity of two sets of data [24, 25]. The dice similarity coefficient is both a spatial overlap indicator and a tool for validating reproducibility. The proportion of specific agreement was another name for it. A DSC value ranges from 0 to 1, with 0 indicating no spatial overlap and 1 representing total overlap. DSC examines the agreement between a predicted segmentation and its ground truth at a pixel level.

The equation for the DSC metric is:

$$DSC = 2 * |X \cap Y| / (|X| + |Y|) \qquad (2)$$

- Where X and Y are two sets.

- A set with vertical bars on either side denotes the set's cardinality, or the number of elements in that set, e.g. |X| denotes the number of elements in set X.

- $\cap$ is used to express the intersection of two sets and refers to the items that are shared by both.

*2) Model Training Optimization Algorithm*

Cost Function

The error between real y and predicted y at its present position is measured by the cost (or loss) function [26]. A loss function can be used to increase the effectiveness of the machine learning model by giving it feedback in order to change the parameters reducing the error, and ultimately locating the local or global minimum. Until the cost function is near to or equal to zero, it iterates repeatedly, moving in the direction of steepest descent. Learning in the model stops at the steepest descent. A cost function determines the average error across the whole training set, whereas a loss function just considers the error of one training sample [27].

Gradient Descent (Gd)

Optimizers update the model in response to the output of the loss function, consequently reducing the loss function. To locate a local minimum or maximum of a given function, gradient descent (GD), an iterative first-order optimization technique, is utilized. This method reduces the cost function and is mostly used as an entry level optimization in machine learning application [28]. Gradient descent however uses a constant learning rate which may need to be tuned manually to reach optimal performance. A higher learning results in faster training which require fewer epochs at the cost of overshooting the minimum. On the contrary a smaller learning rate will result in slower learning which requires more epochs for convergence [29]. The most common optimizers derived from gradient descent are Adaptive Gradient (Adagrad),

Adaptive Delta (AdaDelta), Stochastic Gradient Descent (SGD), Adaptive Momentum (Adam), Cyclic Learning Rate (CLR), Adaptive Max Pooling (Adamax), Root Mean Square Propagation (RMS Prop), Nesterov Adaptive Momentum (Nadam), and Nesterov accelerated gradient (NAG) for CNN [30].

Adaptive Moment Estimation Optimizer (Adam Optimizer)

When training neural networks and machine learning models, the gradient descent approach is typically employed for optimization [27].

The Adam optimizer was made for deep neural network training optimization. It can be described as a combination of momentum-based stochastic gradient descent and RMSprop. Adam optimizer delivers computational performance, lower memory usage, and invariant to diagonal rescaling of gradients for applications with huge amounts of data or parameters [31]. Adam optimizer also computes adaptive learning rate, which entails tuning the learning rate during back-propagation, a property which gives it a competitive edge over other gradient descent optimizers. In study [30], ten common GD based optimizer algorithms were compared and analysed, and results obtained showed that Adam optimizer was more competitive with an accuracy of 99.2%.

*C. Model Conversion to Hardware Description Language (HDL)*

Workflow of implementing neural network in FPGA entails building of model architecture by deciding on the relevant layers and the training of the model using frameworks such as Tensorflow, Keras or Pytorch in high level language such as python and then converts the trained model to hardware description language (HDL). FPGA vendors have high level synthesis tools which synthesize from C++ to HDL. Converting a model in python trained to C++ can be a daunting task. However there are automation tools that bridge between trained models and C++ representations, such as LegUp, DNN Weaver, FINN and HLS4ML among others.

*A. Conversion Tools*

*1) LegUp:* LegUp is a high-level synthesis tool that uses C programming to get the system's behavioural description and creates an RTL netlist in Verilog HDL [32]. LegUp is compatible C/C++ language program and gives an output of Verilog HDL. The four major steps in LegUp HLS process are allocation, scheduling, binding, and RTL generation, which run consecutively [33]. LegUp supports Microchip PolarFire FPGA's.

*2) FINN:* FINN is an experimental framework developed by Xilinix Research Labs that focuses on the use of deep neural networks on FPGAs. It is designed for Xilinx FPGAS for converting quantized neural networks (QNNs) to HDL. FINN supports Brevitas which is a Pytorch package for neural network quantization that supports post training quantization (PTQ) and quantization aware training (QAT) [34].

*3) DNNWeaver:* DNNWeaver is an alternate converter for trained model to HDL implementation in FPGA, made at the

Alternative Computing Technologies (ACT) Laboratory, University of California [35]. A synthesizable FPGA accelerator with high efficiency and performance was achieved in [35]. The tool however only supports Caffe model which is a limitation for Tensorflow and Pytorch frameworks applications.

*4) HLS4ML:* HLS4ML is a tool for machine learning model implementation in FPGA covering both Vivado, Vitis and Quartus HLS backend and C++ inference templates [36]. The tool was developed for the CERN Hadron Collider (LHC), for fast capturing of results from detectors in the LHC. HLS4ML features Keras, Pytorch and ONNX frameworks, models C++ equivalent representations conversion, which can then be transformed to HDL by the backend HLS. Fully connected NN (multilayer perceptron, MLP), Convolutional NN, Recurrent NN (LSTM) and Graph NN (GarNet) are neural network architectures supported by HLS4ML. The support for different backend, frameworks and neural network architectures makes HLS4ML a more suitable tool for FPGA neural network inference.

Neural network implementation consumes a lot of FPGA resources that would be impossible to implement without optimization, which has led research work that incorporates FPGA hardware optimization when using HLS4ML. In study [37], inference of jet substructure model for particle physics in FPGA, archived Flip-flop and BRAM utilization below 4% of budget by employing bit width adjustment. FIFO buffer depth optimization was applied in [38] by recording buffer occupancy during the RTL simulation and then re-running the synthesis with updated FIFO buffer depth reduced resource utilization BRAM by 81%; LUT by 35% and FF by 37%.

## III. METHOD AND MATERIALS

### A. 2D UNET Architecture

The original UNET has five blocks in the contraction path that entail convolution, ReLU and max-pooling operations. Fig. 1 shows four blocks that perform up-sampling, convolution and ReLU in the expansion path, that are connected to their corresponding layers in contracting path. In this work the UNET architecture was reduced to reduce the resource requirement of the model [39].



Fig. 1. UNET architecture [39].

### B. Reduced UNET Model

The depth of the UNET architecture was reduced and the numbers of filters per layer were reduced. Filter numbers were reduced by 93.8% per layer and only one block of filter was used instead of two per layer. In the contraction path the feature extraction blocks were reduced from five to two consisting of 3*3 convolution kernels, 2*2 max-pooling. For the expansion path feature extraction blocks were reduced to two from four consisting of 2*2 up-sampling and concatenate path from the contracting layer as depicted in Fig. 2. The resulting model reduced model total parameters by 99%. Decreasing the model however comes with a cost of accuracy as the features that the model can capture are reduced. Tensorflow Keras framework was used to train the model.



Fig. 2. Reduced UNET architecture.

### C. Data Pre-Processing

Data pre-processing entails formatting dataset images to match the input size and features of the UNET model architecture. The BraTS 2020 dataset was used in this work and was pre-processed by resizing, min-max scaling and slicing before being used for training and testing of the model. In order to capture more features T1CE, T2 and FLAIR modalities were used for training.

Resizing

BraTS 2020 dataset consists of 3D MRI images with modalities and segmentation masks marked by experts. The 3D-MRI scans were resized from 240*240*155 to 128*128*128 dimension to reduce the unnecessary background data and for uniformity in the data. Fig 3 shows the resized image with dimensions 128*128 and background cropped out.

*1) Min-max scaler:* The images were scaled using the min-max scaler for a mean of 0 and max value of 1. Scaling helps to standardize the data for efficient training of the model.

*2) Slicing:* Since the model was defined for 2D convolution kernels each image was sliced to from 3D to 2D slices for all the modalities and masks. The modalities for each sample were converted to numpy array and combined or stacked into one image which translate to input channels in the model architecture. Fig. 3 is a plot of a 2D slice from the 127 slices obtained per 3D image slicing.

Fig. 3.   Pre-processed BraTS 2020 2D slice.

*3) Data Generator:* A data generator was defined for loading the images and masks during training. The images and masks were loaded as an (X,Y) tuple. In the tuple method the X represents the data which was passed through the model for training and Y is the expected output which was used for calculating the loss against the model output at each iteration.

### D. Adam Optimization

Adam optimizer was used during the training to update the weights and learning rate to reach lowest possible loss, through gradient descent method. In order to calculate the loss, dice score coefficient (DSC) was used to compute the total loss which was then back-propagated for updating weights and the learning rate. The optimizer weights were removed from the model after the training since they are not required for inference.

### E. Tensorflow Keras Model Conversion to C++ and RTL

HLS4ML was used for converting the trained model to HDL. Adam optimizer Weights added to the model during training was removed before conversion. This reduced the memory footprint of the model as optimizer weights are not required for inference. A script for generating the C++ equivalence of the model was developed. The script describes the optimization category in terms of resource or latency, clock, input/output type (io_type) and backend HLS tool that was used. Fig. 4 illustrates the process flow from model training to FPGA implementation using HLS4ML.



Fig. 4.   HLS4ML conversion process.

Conversion to RTL

Table I illustrates the configuration used to convert the reduced model to HDL. The C++ generated representation of the model was obtained from the templates in hls4ml that define CONV2D, MAXPOOL and other layers defined in the model. Vivado HLS supports converting C++ model representation to HDL. In this work the C++ model was converted to Verilog, System-C and VHDL. The main files

generated from the conversion are the HDL, data and constraints. HDL files describe various layers of the model; data files contain the weights and biases and the constraint files define the timing and layer interconnections.

TABLE I.        REDUCED UNET CONVERSION TO HDL CONFIGURATION

| Configuration | Parameter |
|---|---|
| Backend | Vivado HLS |
| FPGA | Kintex Ultrascale Part- xcku085-flva1517-3-e |
| Strategy | Resource (Reuse Factor) |
| FIFO Depth (Initial) | 100_000 |
| Precision | Fixed point arithmetic(ap<>) |

### F. HDL Optimization

In order to optimize the FPGA resource utilization the precision, resource strategy and FIFO buffer depth optimization techniques were used.

*1) Precision:* Floating point numbers due to their limitless precision in computation leads to an increased utilization of resources. Arbitrary fixed point type was used in this work which is defined by 'ap_fixed<a,b>', where 'a' is the total bit width and 'b' is the fractional part from the total size of 'b'. Fixed-point arithmetic is more efficient in reducing resource utilization when compared to floating-point arithmetic.

*2) Resource strategy:* In multilayer neural networks, each neuron in a layer (consisting of n neurons) produces an output computed by the weighted sum of the output from the previous layer. The weights associated with these connections are represented as a matrix W of size n*m, where m is the number of neurons in the previous layer. Each neuron has a bias that is independent and represented by vector b. The weighted sums and the biases are summed and added non linearity by activation function denoted by f, resulting in the final output of the neuron. This process of weights, biases and activation function is expressed as:

$$y=f(W*\text{previous layer output} + b) \qquad (3)$$

The multiplications in neural networks are computed by multipliers in FPGA, which result in high resource utilization if each multiplication was to represent a physical multiplier. Resource strategy was used to optimize the design by reusing of multipliers as demonstrated in Fig. 5.



Fig. 5.   Multiplier reuse reduces number of multipliers at the expense of parallel processing.

Multiplier Limit

Processing elements are a functional block that preforms specific operations such as multiplication and addition. In the FPGA, operations are efficiently conducted in parallel across multiple PE's, thereby improving computational efficiency. For resource strategy multiplier_limit variable was defined to represent the maximum number of multiplications that can be done in parallel for the available resources. The calculation is based on the number of input values (mult_n_in), the number of output values ( mult_n_out) and the reuse factor. A higher reuse factor will mean a lower value of multiplier_limit which translates to actual multipliers used in the FPGA implementation, however the universal time of operation will increase as the multipliers will be shared by a number of operations.

$$Mult\_n\_in = filt\_height * filt\_width * n\_chan \quad (4)$$

$$Mult\_n\_out = n\_filt \quad (5)$$

$$Multiplier\_limit = DIV\_ROUNDUP(mult\_n\_in * mult\_n\_out, reuse\_factor) \quad (6)$$

*3) First In First Out (FIFO) Buffer Depth Optimization:* FIFO buffers store data in between layers. In the initial C++ conversion the size of the buffers is estimated, however the estimated size is above the utilization during simulation. This results in a BRAM and LUT usage that is higher than what the design requires. In FIFO depth optimization method the buffer size for each layer was set to 100_000, which is a value above what is required by the design. The design was then synthesized and during RTL co-simulation the buffer occupation was recorded and used to update the FIFO buffer sizes which translated to a 71% reduction in buffer size. Fig. 6 is an overview of the buffer method capturing of buffer utilization during simulation to the updating of the new buffer size.



Fig. 6. FIFO buffer depth optimization overview.

*G. C++ /Register Transfer Level (C/RTL) CO-SIMULATION*

C/RTL co-simulation was used to verify the functional preservation of the HDL converted model. The test images were converted to a 1D array and saved as data files with pixel values saved as strings for compatibility with simulator. C/RTL co-simulation was done using test-bench and test data. The comparison of C-Simulation and RTL simulation passed

validation and output for inference was a 1D array that was converted to 2D numpy array.

## IV. RESULTS

*A. Precision*

The increase in configured precision for the reduced UNET model during conversion to HDL was directly proportional to the increase in resource utilization as shown in Table II.

TABLE II. INCREASING AP PRECISION INCREASES RESOURCE UTILIZATION FOR THE REDUCED UNET MODEL CONVERSION IN HLS4ML

| PRECISION | BRAM_18K | DSP48E | FF | LUT |
|---|---|---|---|---|
| <16,6> | 2118 | 5 | 67033 | 149764 |
| <32,6> | 3978 | 13 | 97408 | 188001 |
| <64,6> | 7845 | 65 | 155388 | 219201 |

*B. Resource Strategy*

Multiplier reuse reduced the resources used as multipliers were shared among operations at the expense of performance. Maximum possible reuse factor is 1152 hence beyond 1500 reuse factor the performance and resource utilization is not affected, as shown by Fig. 7 graph of resource utilization. This is because there is an upper-limit to optimization of resource utilization against performance.



Fig. 7. Increase in reuse factor is directly proportional to resource utilization and reach roofline at 1000 reuse factor.

*C. First In First Out (FIFO) Buffer Depth Optimization*

New FIFO buffer depth values were inserted in the C++ firmware as pragma directives to guide synthesis. Resource utilization after FIFO optimization showed 54.8% reduction in BRAM's, 29% reduction in FF's and 44% reduction in LUT's. The new FIFO buffer depth was the occupancy increased by value of 1. Table III shows the reduction in resource utilization pre and post buffer size reduction. Table IV illustrates how the buffer size was reduced using pragma directives during synthesis.

TABLE III. FIFO DEPTH OPTIMIZATION REDUCTION IN RESOURCES

| Optimization | BRAM_18K | DSP48E | FF | LUT | LATENCY(ms) |
|---|---|---|---|---|---|
| No | 1634 | 8 | 52679 | 120583 | 9.8 |
| Yes | 739 | 5 | 29811 | 62059 | 9.8 |
| % Change | -54.8 | -37.5 | -43.4 | -48.5 | 0 |

TABLE IV.    FIFO BUFFER DEPTH SIMULATION OCCUPANCY RESULTS AND OPTIMIZATION PRAGMAS TO REDUCE BUFFER DEPTH

| Layer | Occupancy | New FIFO BUFFER DEPTH | PRAGMA DIRECTIVE |
|---|---|---|---|
| layer20_out_V_data_0_V_U | 33348 | 33349 | #pragma HLS STREAM variable=layer20_out depth=33349 |
| layer21_out_V_data_0_V_U | 16608 | 16609 | #pragma HLS STREAM variable=layer21_out depth=16609 |
| layer19_cpy2_V_data_0_V_U | 819 | 820 | #pragma HLS STREAM variable=layer19_cpy2 depth=820 |
| layer22_out_V_data_0_V_U | 134 | 135 | #pragma HLS STREAM variable=layer22_out depth=135 |
| layer23_out_V_data_0_V_U | 199 | 200 | #pragma HLS STREAM variable=layer23_out depth=200 |
| layer12_out_V_data_0_V_U | 3 | 4 | #pragma HLS STREAM variable=layer12_out depth=4 |
| layer24_out_V_data_0_V_U | 37372 | 37373 | #pragma HLS STREAM variable=layer24_out depth=37373 |
| layer25_out_V_data_0_V_U | 261 | 262 | #pragma HLS STREAM variable=layer25_out depth=262 |
| layer26_out_V_data_0_V_U | 1 | 2 | #pragma HLS STREAM variable=layer26_out depth=2 |

TABLE V.    SHOWING THE REDUCED MODEL COMPARISON WITH EXISTING WORK

| Ref | Model | Dataset | Parameters | IoU Score | Accuracy | Execution Time(ms) |
|---|---|---|---|---|---|---|
| UNET MODEL | UNET | BraTS 2020 | 1.9 M | 0.60 | 0.97 | 132 |
| [40]Ercüment GÜVENÇ (2023) | FLAIR MR IMAGES WITH U-NET | BraTS 2018 | - | 0.59 | 0.99 | - |
| [41]Jwaid (2021) | 3D U-Net CNN | BraTS 2017 | - | 0.69 | 0.99 | - |
| Proposed | Reduced UNET | BraTS 2020 | 864 | 0.74 | 0.95 | 38 |

### D. Reduced Model Result Comparison with Existing Work

A full UNET model with original architecture was developed to compare with the reduced UNET proposed in this work. Model parameters were reduced from 1.9 million in original 2D-UNET model to 864 in proposed reduced UNET model.

### E. Intellectual Property (IP) Core

Fig. 8 is the graphical representation of the intellectual property core that was generated from the RTL design. The generate IP core has 3-inputs which correspond to the three input channels in the model architecture and there are four output channels which are aligned to the number of classes for type of tumour and background. IP cores introduce modular design which can then be interfaced with input and output IP cores such as Ethernet.



Fig. 8.    BraTS 2D UNET tumour segmentation IP core.

### F. C/RTL Co-simulation Results

The co-simulation results presented in Fig. 9, visually demonstrate the similarity between prediction output masks of the python model and RTL simulation. Quantitative analysis, computing the IoU score output a comparable value of 74% between the original python model and the RTL simulation, validating the fidelity of the FPGA-based implementation.

Fig. 9.    Reduced UNET python model and C/RTL co-simulation brain tumour prediction masks.

## V. DISCUSSION

Table V shows that the proposed model has a higher IoU score over the existing work in [40, 41]. In study [40], only FLAIR images modality which limited the amount of features that can be learned by the model as compared to T1CE, T2 and FLAIR modalities which were used in the proposed model to learn more features from the dataset. The proposed model however has limitations in terms of accurately predicting tumour class. Future work can be focused on improving multiclass prediction of the proposed model.

## VI.    CONCLUSION

In this work a reduced UNET model was built and trained in Tensorflow Keras for brain tumour segmentation applications and archived an IoU score of 74%. The reduced model significantly reduced the model parameters by 99% which translates into reduced computational requirements. Converting the reduced model into C++ and HDL equivalent representations using HLS4ML, FPGA resources were used economically while preserving the original model segmentation output mask accuracy. Resource strategy, FIFO buffer depth and precision methods significantly reduced FPGA resource usage. The reduced model segmentation performs well in predicting tumour region, however the tumour classes are still poorly predicted. Further work on choosing the right loss function suitable for unbalanced multi-class segmentation for the reduced model can be done to improve inference accuracy while reducing FPGA resource utilization

## REFERENCES

[1]    Clinic, Mayo, "Brain tumor," Mayo Clinic, 2023. [Online]. Available: https://www.mayoclinic.org/diseases-conditions/brain-tumor/symptoms-causes/syc-20350084. [Accessed 20 November 2023].

[2]    S. R. H. K. S. R. N. K. Soheila Saeedi, "MRI-based brain tumor detection using convolutional deep learning methods and chosen machine learning techniques," BMC Medical Informatics and Decision Making, vol. 23, no. 16, 2023.

[3]    Cancer.Net, "Brain Tumor: Diagnosis," Cancer.Net, March 2023. [Online]. Available: https://www.cancer.net/cancer-types/brain-tumor/diagnosis. [Accessed 21 November 2023].

[4]    A. Abdullah Asiri et al, "Brain Tumor Detection and Classification Using Fine-Tuned CNN with ResNet50 and U-Net Model: A Study on TCGA-LGG and TCIA Dataset for MRI Applications," Artificial Intelligence Applications in Medical Imaging, vol. 13, no. 7, 2023.

[5]    A. K. T. M. P. S. A. B. Mukul Aggarwal, "An early detection and segmentation of Brain Tumor using Deep Neural Network," BMC Medical Informatics and Decision Making, vol. 23, no. 78, 2023.

[6]    P. F. T. B. Olaf Ronneberger, "U-Net: Convolutional Networks for Biomedical," in International Conference on Medical Image Computing and Computer-Assisted Intervention, 2015.

[7]    Xuan Cheng et al, "Efficient hardware design of a deep U-net model for pixel-level ECG classification in healthcare device," Microelectronics Journal, vol. 126, 2022.

[8]    Z. L. Chenghao Wang, "A Review of the Optimal Design of Neural Networks Based on FPGA," MDPI, p. 44, 2022.

[9]    M. B. B. K. A. K. Dinthisrang Daimary, "Brain Tumor Segmentation from MRI Images using Hybrid Convolutional Neural Networks," in International Conference on Computational Intelligence and Data Science (ICCIDS 2019), 2020.

[10]   E. B. V. D. C. Tonya White, "Data sharing and privacy issues in neuroimaging research: Opportunities, obstacles, challenges, and monsters under the bed," Hum Brain Mapping, vol. 43, no. 1, pp. 278-291, 2020.

[11]   K. T. K. A. Keiko Sato, "Five Genes Associated With Survival in Patients With Lower-grade Gliomas Were Identified by Information-theoretical Analysis," Anticancer, pp. 2777-2785, 2020.

[12]   Emmanuel Rios Velazquez et al, "Fully automatic GBM segmentation in the TCGA-GBM dataset: Prognosis and correlation with VASARI features," Scientific Reports, vol. 5, no. 16822, 2015.

[13]   P. C. W. Katarzyna Tomczak, "The Cancer Genome Atlas (TCGA): an immeasurable source of knowledge," Contemporary Oncology, vol. 19, no. 1A/2015, pp. A68-A77, 2015.

[14]   Spyridon Bakas et al, "Advancing The Cancer Genome Atlas glioma MRI collections with expert segmentation labels and radiomic features," Scientific Data, vol. 4, no. 1, 2017.

[15]   BrainWeb, "BrainWeb: Simulated Brain Database," McGil University, [Online]. Available: https://brainweb.bic.mni.mcgill.ca/. [Accessed 22 May 2023].

[16]   A. Z. Elnomery, "An Adaptive Fuzzy C-Means Algorithm for Improving MRI Segmentation," Open Journal of Medical Imaging, vol. 3, pp. 125-135, 2013.

[17]   S. Bakas, "Multimodal Brain Tumor Segmentation Challenge 2020: Data," Perelman School of Medicine, University of Pennsylvania, 2020. [Online]. Available: https://www.med.upenn.edu/cbica/brats2020/data.html. [Accessed 22 May 2023].

[18]   Bjoern H. Menze et al, "The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS)," IEEE Transactions on Medical Imaging, vol. 34, no. 10, pp. 1993-2024, 2015.

[19]   Muhammad Usman Saeed et al, "RMU-Net: A Novel Residual Mobile U-Net Model for Brain Tumor Segmentation from MR Images," Electronics, vol. 10, no. 16, 2021.

[20]   Spyridon Bakas et al, "Identifying the Best Machine Learning Algorithms for Brain Tumor Segmentation, Progression Assessment, and Overall Survival Prediction in the BRATS Challenge," arXiv preprint arXiv:1811.02629, 2019.

[21]   S. U. A. V. J. R. Arumuga Arun, "Reduced U-Net Architecture for Classifying Crop and Weed using Pixel-wise Segemntation," in 2020 IEEE International Conference for Innovation in Technology (INOCON), Bengaluru, India, 2020.

[22]   M. S. R. Rupal R. Agravat, "Prediction of Overall Survival of Brain Tumor Patients," in TENCON, Kochi India, 2019.

[23] O. Sheremet, "Intersection over union (IoU) calculation for evaluating an image segmentation model," Towards Data Science, 25 July 2020. [Online]. Available: https://towardsdatascience.com/intersection-over-union-iou-calculation-for-evaluating-an-image-segmentation-model-8b22e2e84686. [Accessed 02 July 2023].

[24] B. D. J, "Dice similarity coefficient," Radiopaedia, 02 08 2021. [Online]. Available: https://radiopaedia.org/articles/dice-similarity-coefficient. [Accessed 10 02 2023].

[25] D. L. C. I. P. F. K. Anthony D. Yao, "Deep Learning in Neuroradiology: A Systematic Review of Current Algorithms and Approaches for the New Wave of Imaging Technology," Radiology: Artificial Intelligence, vol. 2, no. 2, pp. 1-6, 2020.

[26] D. S. S. L. X. L. Yingjie Tian, "Recent advances on loss functions in deep learning for computer vision," Neurocomputing, vol. 497, pp. 129-158, 2022.

[27] IBM, "What is gradient descent?," IBM, 2023. [Online]. Available: https://www.ibm.com/topics/gradient-descent. [Accessed 01 July 2023].

[28] R. Kwiatkowski, "Gradient Descent Algorithm — a deep dive," Towards Data Science, 22 May 2021. [Online]. Available: https://towardsdatascience.com/gradient-descent-algorithm-a-deep-dive-cf04e8115f21. [Accessed 01 July 2023].

[29] J. Brownlee, "Understand the Impact of Learning Rate on Neural Network Performance," MACHINE LEARNING MASTERY, 25 January 2019. [Online]. Available: https://machinelearningmastery.com/understand-the-dynamics-of-learning-rate-on-deep-learning-neural-networks/. [Accessed 01 July 2023].

[30] Muhammad Yaqub et al, "State-of-the-Art CNN Optimizer for Brain Tumor Segmentation in Magnetic Resonance Images," Brain Sciences, vol. 10, no. 7, p. 427, 2020.

[31] J. L. B. Dlederlk P.Kingma, "ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION," in 3rd International Conference for Learning Representations, San Diego, 2015.

[32] S. R. a. M. Joseph, "Open source HLS tools: A stepping stone for modern electronic CAD," in 2016 IEEE International Conference on Computational Intelligence and Computing Research (ICCIC), Chennai, India, 2016.

[33] Inc, Microchip Technology, "2User Guide: 2.1 Introduction to High-Level Synthesis," Microchip Technology Inc, 2021. [Online]. Available: https://download-soc.microsemi.com/FPGA/HLS-EAP/docs/legup-2021.1-docs/userguide.html. [Accessed 30 May 2023].

[34] pypi, "Quantization-aware training in PyTorch," pypi, 28 April 2023. [Online]. Available: https://pypi.org/project/brevitas/#:~:text=Brevitas%20is%20a%20PyTorch%20library,not%20an%20official%20Xilinx%20product.. [Accessed 27 May 2023].

[35] Hardick Sharma et al, "From High-Level Deep Neural Models to FPGAs," in The 49th Annual IEEE/ACM International Symposium on Microarchitecture, Taipei, Taiwan, 2016.

[36] Team, FastML, "fastmachinelearning/hls4ml," FastML Team, 2023. [Online]. Available: https://github.com/fastmachinelearning/hls4ml. [Accessed 27 October 2023].

[37] Javier Duartea et al, "Fast inference of deep neural networks in FPGAs for particle physics," Journal of Instrumentation, vol. 13, no. 07, p. P07027, 2018.

[38] Nicolò Ghielmetti et al, "REAL-TIME SEMANTIC SEGMENTATION ON FPGAS FOR AUTONOMOUS VEHICLES WITH HLS4ML," Machine Learning: Science and Technology, vol. 3, no. 4, 2022.

[39] J. Zhang, "UNET- Line by Line Explanation," Towards Data Science, 4 October 2019. [Online]. Available: https://towardsdatascience.com/unet-line-by-line-explanation-9b191c76baf5. [Accessed 05 November 2023].

[40] M. E. G. Ç. Ercüment GÜVENÇ, "BRAIN TUMOR SEGMENTATION ON FLAIR MR IMAGES WITH U-NET," Mugla Journal of Science and Technology, vol. 9, no. 1, pp. 34-41, 2023.

[41] W. M. Jwaid et al, "Development of Brain Tumor Segmentation of," Eastern-European Journal of, vol. 4, no. 9, pp. 23-31, 2021.

# Enhancing Diabetes Management: A Hybrid Adaptive Machine Learning Approach for Intelligent Patient Monitoring in e-Health Systems

Sushil Dohare[1], Dr.Deeba K[2], Laxmi Pamulaparthy[3], Shokhjakhon Abdufattokhov[4],
Janjhyam Venkata Naga Ramesh[5], Prof. Ts. Dr. Yousef A.Baker El-Ebiary[6], Dr.E.Thenmozhi[7]
Department of Epidemiology, College of Public Health and Tropical Medicine, Jazan University, Saudi Arabia[1]
Assistant Professor, School of Computer Science and Applications, REVA University, Bangalore[2]
Assistant Professor, Vignana Bharathi Institute of Technology, Ghatkesar, Hyderabad.
Affiliated to JNTUH (Autonomous)[3]
Automatic Control and Computer Engineering Department, Turin Polytechnic University in Tashkent, Tashkent, Uzbekistan,
Department of Information Technologies, Tashkent International University of Education, Tashkent, Uzbekistan[4]
Assistant Professor, Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation,
Vaddeswaram, Guntur Dist., Andhra Pradesh - 522302, India[5]
Faculty of Informatics and Computing, UniSZA University, Malaysia[6]
Associate Professor, Department of Information Technology, Panimalar Engineering College, Chennai, India[7]

*Abstract*—The goal of the present research is to better understand the need of accurate and ongoing monitoring in the complicated chronic metabolic disease known as diabetes. With the integration of an intelligent system utilising a hybrid adaptive machine learning classifier, the suggested method presents a novel way to tracking individuals with diabetes. The system uses cutting edge technologies like intelligent tracking and machine learning (ML) to improve the efficacy and accuracy of diabetes patient monitoring. Integrating smart gadgets, sensors, and telephones in key locations to gather full body dimension data that is essential for diabetic health forms the architectural basis. Using a dataset that includes comprehensive data on the patient's characteristics and glucose levels, this investigation looks at sixty-two diabetic patients who were followed up on a daily basis for sixty-seven days. The study presents a hybrid architecture that combines a Convolutional Neural Network (CNN) with a Support Vector Machine (SVM) in order to optimise system performance. To train and optimise the hybrid model, Grey Wolf Optimisation (GWO) is utilised, drawing inspiration from collaborative optimisation in wolf packs. Thorough assessment, utilising standardised performance criteria including recall, F1-Score, accuracy, precision, and the Receiver Operating Characteristic (ROC) Curve, methodically verifies the suggested solution. The results reveal a remarkable 99.6% accuracy rate, which shows a considerable increase throughout training epochs. The CNN-SVM hybrid model achieves a classification accuracy advantage of around 4.15% over traditional techniques such as SVM, Decision Trees, and Sequential Minimal Optimisation. Python software is used to implement the suggested CNN-SVM technique. This research advances e-health systems by presenting a novel framework for effective diabetic patient monitoring that integrates machine learning, intelligent tracking, and optimisation techniques. The results point to a great deal of promise for the proposed method in the field of medicine, especially in the accurate diagnosis and follow-up of diabetic patients, which would provide opportunities for tailored and adaptable patient care.

*Keywords*—*Diabetes; machine learning; convolutional neural network; support vector machine; grey wolf optimization; e-health systems*

## I. INTRODUCTION

Over the past few decades, diabetes mellitus, frequently characterized to as diabetes, has become a major worldwide health concern due to its constantly rising predominance. Increased levels of glucose in the blood are a characteristic of this metabolic illness, which is brought on by either inadequate insulin synthesis or an inefficient utilization of insulin by the body. The World Health Organization (WHO) reports, that the rate of diabetes has been rapidly rising globally, making it one of the biggest public health issues of the twenty-first century. Diabetes affects a wide range of individuals worldwide, as evidenced by its epidemiology [1]. With approximately 463 million individuals identified with diabetes as of 2019, developed nations as well as developing ones are struggling with an increase in the number of instances of the disease. If current conditions continue, this number is expected to rise, reaching an astounding 700 million people by 2045. Inactive ways of life, inadequate nutrition, and a growing elderly population are all contributing participants to this trend, which emphasizes the critical requirement for efficient management and preventive actions [2]. Diabetes that goes untreated has serious repercussions that impact several organ systems and cause significant health issues. Long-term effects include renal failure, blindness, neuropathy, and cardiovascular disorders, all of which significantly decrease the standards of life as well as the life expectancies of those who are impacted. Furthermore, the financial impact of diabetes-related medical expenses and lost productivity is significant, creating new difficulties for healthcare systems throughout the world [3].

There are several interrelated causes that lead to the diabetes pandemic. With the introduction of diets excessive in

fats that are unhealthy and sugars that are processed and a decrease in physical activity, industrialization and modifications to lifestyles have created an environment that encourages obesity that promotes the development of diabetes. An individual's risk to the illness is influenced by both environmental and genetic variables, which combine to determine the person's susceptibility. Furthermore, the growing incidence of diabetes is made worsened by differences in the availability of healthcare and education, especially for communities with limited resources. Diabetes is becoming more and more common, which has serious consequences for public health and necessitates an all-encompassing strategy [4]. The three main strategies for reducing the rise in diabetes are prevention via health education, lifestyle modifications, as well as early detection. In addition, technological developments including the development of sophisticated monitoring systems have the potential to improve diabetes care while promoting an anticipatory approach to healthcare [5]. To reduce the effects of diabetes and enhance the physical and mental well-being of millions of people globally, the international community needs to collaborate together to address the complex interactions between hereditary, environmental, and lifestyle variables [6].

The complex nature of managing diabetes, characterized by the requirement for continual surveillance and individualized treatment, highlights the necessity for sophisticated monitoring systems. Diabetes, in contrast to many other chronic illnesses, requires close monitoring of blood sugar levels, dietary habits, activity levels, and medication compliance. Diabetes has several facets that impose significant stress on patients and healthcare professionals equally [7]. As a result, there is a strong demand for innovative approaches that may expedite monitoring procedures, give real-time information, and enable faster treatments. Diabetes is a very unique disorder, with changes in medical condition and patient responses to therapy occurring on an individual basis. In order to customize treatment plans to the particular requirements of each patient, sophisticated monitoring systems are now essential [8]. Through the integration of adaptive learning technologies and sophisticated monitoring mechanisms, these systems are able to evaluate large datasets and identify specific trends, offering a more detailed knowledge of a patient's health trajectories. Enhancing therapeutic efficacy, reducing adverse reactions, and ultimately improving patient outcomes are all possible with this customized strategy [9].

Patient-centered care is an innovative approach that prioritizes giving individuals the tools required to take an active role in their own health management. Modern monitoring devices are essential to this change because they provide patients with immediate information on their lifestyle decisions and health parameters. These systems have the potential for motivating patients to follow treatment programs, make educated decisions, and establish up healthy habits by cultivating an environment of ownership and awareness [10]. Furthermore, the incorporation of easy-to-use interfaces and smartphone applications might promote a proactive and cooperative approach to diabetes management by facilitating effortless interaction between patients and medical

professionals. By using sophisticated technologies to monitor diabetes proactively, complications may be avoided, and the financial strain of the disease may be reduced. Early intervention can be used to mitigate the frequency of serious illnesses and hospitalizations by promptly detecting abnormalities from normal health indicators [11]. Sophisticated monitoring systems help healthcare systems preserve revenue over the course of time through promoting preventative care and supporting continuous maintenance of health. In order to effectively manage the numerous obstacles presented by this complicated and widespread chronic illness, improved monitoring systems are becoming increasingly necessary as the number of cases of diabetes rises worldwide [12].

Significant progress has been made in the field of diabetes e-health systems, which use technology to improve the treatment of patients and management. These systems usually incorporate a range of technologies to monitor and assist people with diabetes in everyday activities, such as online platforms, wearable technology, and mobile applications. Numerous current systems concentrate on monitoring blood glucose levels, activity levels, and consumption habits in real-time. While wearable technology, such continuous glucose monitors (CGMs), offer a constant supply of physiological information, mobile applications frequently act as a central centre for data gathering and processing [13]. e-Health technologies have made it possible for diabetes patients to take advantage of telehealth services and remote monitoring, eliminating the distance between patients and healthcare professionals in geographically dispersed areas. Regular assessments and treatment plan modifications are made possible through telehealth conversations, which eliminate the need for several visits to the clinic. By facilitating more adaptable and dynamic diabetic treatment, this integration lowers the need for frequent visits to the clinic and increases patient participation.

The incompatibility of various e-health technologies and systems is one of the main disadvantages. Variations in information formats and requirements may prevent the effortless transfer of information between healthcare systems and devices, which could result in evaluations of an individual's health state that are either inaccurate or incomplete. Security and confidentiality of information are major problems in e-health systems because of the sensitive character of health data. Challenges including illegal entry, leaks of information, and insufficient encryption protections caused patient privacy at risk and might undermine users' confidence in these systems. Participation among users and commitment over time to monitoring methods remain issues in regardless of the capacities of e-health systems [14]. The apparent complexities of the systems, unease with wearable technology, or a lack of customized input that aligns with their specific health objectives can all lead to patients detaching from ongoing surveillance. Certain e-health systems could depend on standard algorithms that do not adequately take into consideration the range of diabetes appearances people with the disease can have. An approach that applies to all patients may fail to recognize smaller variations in how each patient reacts to therapy, which could result in undesirable outcomes

for certain groups of patients. e-Health systems are effective at collecting physiological data, but they can still do better when it comes to incorporating behavioural data, including how anxiety or other lifestyle factors affect diabetic treatment [15]. Increasing the comprehension of every facet of the patient experience can result in more comprehensive and individualized treatments. Improving interoperability, enhancing privacy and security protocols, boosting user engagement with intuitive interfaces, and fine-tuning algorithms in order to accommodate a range of patient profiles are all necessary to overcome these constraints. The continued development of e-health systems shows possibilities for improving patient outcomes, optimizing diabetic treatment, and expanding the field of electronic health records as technology advances.

The complicated and constantly changing nature of diabetes is the motivating factor driving the integration of intelligent tracking and a hybrid adaptive machine learning classifier in diabetes treatment. Individualized patterns of blood glucose levels, choices regarding lifestyle, and treatment responses are characteristics of diabetes. The complex and changing health trajectories of individuals with diabetes are frequently difficult for traditional, static models to represent. The system attempts to give patients a more comprehensive and individualized approach to diabetes management by combining intelligent tracking, which continually records and customizes to changing patient behaviours, and a hybrid adaptive machine learning classifier, which can learn complicated patterns in large datasets [16]. The traditional standardized method of managing diabetes may not be able to adequately satisfy each patient's specific demands. The device can collect instantaneous information on a patient's activities, eating habits, and physiological reactions because of intelligent tracking systems. This is enhanced by the hybrid adaptive machine learning classifier, which gains knowledge from the recorded information, customizes its model to account for individual differences, and offers individualized recommendations. This combination maximizes the effectiveness of therapies and improves patient outcomes by facilitating the transition towards more individualized and focused treatment techniques.

Continuous monitoring helps significantly in the treatment of diabetes since it makes it possible to identify small variations in health indicators that might signal impending emergencies. Intelligent tracking is integrated to provide an uninterrupted supply of pertinent data, and the hybrid adaptive machine learning classifier is extremely skilled at identifying complex patterns linked to early indicators of health decline. For those with diabetes, this continuous surveillance and early intervention strategy may help avoid complications, lessen the need for emergency interventions, and enhance their general quality of life. The hybrid adaptive machine learning classifier is intended to address the difficulties caused by the intrinsic unpredictability in the initial responses of individuals with diabetes to dietary and medication modifications [17]. Conventional classifiers could have trouble adjusting to these differences, which would result in less than ideal efficiency. Because of its adaptive characteristics, the suggested classifier can adapt over time to the intrinsic variety in diabetes

presentations. This flexibility is especially important for patients with chronic conditions like diabetes, whose health can be affected by a wide range of variables. In the area of e-health and diabetes care, the combination of intelligent tracking with a hybrid adaptive machine learning classifier is a newly developed and creative method. While discrete components like machine learning and intelligent tracking have been studied independently, integrating them into a unified framework which functions effectively when combined is a novel contribution. The system's capacity to dynamically adjust to each patient's unique profile, learn from changing information over time, and offer customized suggestions for successful diabetes control essentially makes this system exceptional. This strategy might contribute to the expanding area of customized medicine and raise the standards for digital health interventions for chronic illnesses.

The Key Contribution of the paper is given as follows:

- The research presents a unique hybrid architecture that combines a SVM with a CNN to monitor diabetic patients. By combining the attributes of both models which is SVM's outstanding binary categorization and CNN's feature extraction capabilities, improves diabetes prediction accuracy.

- An intelligent tracking mechanism is utilized to collects detailed body dimension information from diabetes patients using cellphones, sensors, and smart devices. Beyond typical monitoring techniques, this integrated strategy ensures a more comprehensive awareness of the patient's health and contributes to a more individualized approach to treatment.

- The employed Grey Wolf Optimization is based on cooperative optimization observed among wolf packs, to fine-tune the hybrid model. This optimization method improves the efficacy and effectiveness of the model, offering an innovative approach for fine-tuning parameters in machine learning systems that are inspired by nature.

- The suggested approach shows a significant increase in accuracy. CNN-SVM hybrid model exhibits improved classification accuracy when compared to established approaches such as SVM, Decision Trees, and Sequential Minimal Optimization. This indicates the model's potential for dependable diabetic patient monitoring.

- Machine Learning, optimization, and intelligent tracking approaches, contributed to the improvement of e-health systems. This strategy has a lot of potential to improve the accuracy and efficiency of diabetes patient monitoring. Modern technology combined with the enhanced performance of the suggested hybrid model is a significant addition to the field of healthcare informatics.

The rest of the section is organised as shown below. Section II illustrates literature works on e-health Systems. Section III gives the Problem Statement. Section IV covers the proposed technique for Monitoring and Tracking the Patients

with Diabetes. Section V illustrates the performance measures and summarises the findings and compares the method's performance to previous techniques. Section VI summarises the conclusion and paves the way for future works.

## II. RELATED WORKS

Individuals with diabetes who receive continuous medical attention often have a greater standard lifestyle than those who do not. Due to technology improvements, healthcare costs can be reduced by utilizing the Internet of Things. Both the advancement of intelligent devices and a growth in the total amount of software linked to the networks are necessary for addressing the demands of e-health applications. Therefore, the cellular network must be able to accommodate sophisticated medical applications which require outstanding energy consumption in order to accomplish these objectives. The study develops combined voting classifier that utilizes neural networks to effectively forecast patients' diabetes through online monitoring. Internet of Things gadgets is used in the study to track patient cases. IoT devices provide their information for smartphones during evaluating, and those devices transmit the information to the cloud, where categorization is done. The Python tool is used to run the simulation on the gathered observations. The results from the simulation demonstrate that, in comparison to current the most advanced combination models, the suggested strategy provides a higher prediction rate, precision, recall, and f-measure. However, for information to be provided from connected devices to the cloud, the suggested solution depends on the Internet operating without interruptions. Restrictions in internet access might possibly undermine the dependability of the projections by affecting the continuous monitoring capability [18].

Because type 2 diabetes has a significant condition of disease and significantly lowers the patient's standard of existence, using computerized instruments and data technologies to control illness has become common place due to the close connection between healthcare and the worldwide web. The study attempted to determine how well several e-health treatments, varying in length, may help individuals who have type 2 diabetes achieve controlling their glycemic levels. Researchers investigated for randomization managed studies describing various e-health interventions for glycemic management in individuals with type 2 diabetes. Individuals with type 2 diabetes mellitus satisfied the following participation requirements: (1) interventional duration ≥1 month; (2) findings HbA1c (%); and (4) randomization management using e-health based techniques. Cochrane techniques were employed to evaluate potential biases. Researchers performed the Bayesian network a meta-examination using R 4.1.2. The most efficient intervention periods were found to be ≤6 months, according to subgroup evaluations. Individuals with diabetes who have type 2 diabetes can benefit from improved glycemic management through all forms of e-health-based interventions. With an ideal intervention length of ≤6 months, SMS is a high-frequency signal, low-barrier technique that delivers the highest benefit in decreasing HbA1c. The research concepts administration procedures and characteristics of participants may add variability due to the numerous types and timings of e-health treatments included in the inclusion. The variety of approaches may make it more difficult to reach firm judgments on the efficacy of particular therapies [19].

Many of the advances that technology has enabled about for humankind have made even the most difficult things simpler. The development of science and technology has led to the widespread deployment of intelligent machines. Numerous sectors, including health care for the public, medicine, and healthcare, have seen advancements. The monitoring of wellness and various other tasks may now be done effectively, economically, and intelligently thanks to recent advancements in healthcare. This is made feasible in large part by wearables. Despite their compact design, these gadgets are equipped with potent medical sensors that enable the monitoring of users' health problems. As technology and medical science have advanced, wearables have been equipped with an increasing number of sensors for tracking a wide range of activities, such as blood oxygen levels, body temperatures, cardiovascular disease, and activity monitoring systems, among many others. These gadgets allow users to store the outcomes for future usage and can be linked to smartphones. The necessity of wearables in healthcare and their potential to transform medical systems in the future were covered in this study. A case study is given illustrating how wearables may benefit both physicians and patients. Potential applications and research problems are also included in this publication. However, because patient groups vary and wearable devices are different, the report may have difficulty generalizing its results. The general application of the offered findings is limited by the realization that wearable device performance and utilization could differ across various populations, medical conditions, and geographical areas [20].

Globally, an advancing population is one of the biggest healthcare challenges. Due to their increased risk of chronic illnesses, which raise healthcare costs, older individuals demand greater resources from the healthcare system. One of the major developments within healthcare technological advances is the creation and ongoing operation of e-health solutions, which provide patients with mobile services to support and improve their treatment according to monitoring certain physiological information. Healthcare technology has advanced greatly in the previous few decades in terms of size, velocity, accessibility, and connectivity. The fact that individuals are restricted to smart rooms and a bed equipped with monitoring devices is a significant disadvantage of contemporary e-health monitoring systems. Because chronic patients have accessibility, security, and adaptability difficulties, such surveillance is not common. Furthermore, attached to a patient's body health monitoring gadgets provide no evaluations or recommendations. This work presents a multi-agent-based approach to tracking health conditions that aims to enhance the procedure by gathering patient information, reasoning collectively, and suggesting actions to patients as well as physicians in a mobile setting. The paper presents a multi-agent-based system that is assessed using a case study. The findings demonstrate that the suggested approach offers elderly, chronic, and distant patients an effective way to monitor their health. Furthermore, by employing 5G technology, the suggested method works better

than the current health system and enables prompt medical services for patients who live far away. However, when the suggested agent-based system is expanded to support enormous user bases or extensive medical networks, its efficacy can encounter difficulties. As the size grows, problems with resource use, communication costs, and system efficiency could become increasingly noticeable [21].

Diabetes is a chronic illness caused by the pancreas' inability to produce enough insulin or to shield the body from non-consumable substances. Diabetes patient health monitoring is a methodical approach that provides us with comprehensive health information on individuals with diabetes. Health observation platforms for diabetic patients are essential for monitoring their state, especially when using Internet of Things connected devices. In simple terms, diabetic patient monitoring platforms may screen individuals with diabetes and save certain health data, such as body temperature, blood pressure, and blood glucose levels. Because predictive analysis may assist diabetic individuals, their relatives, medical professionals, and clinical researchers in making decisions about the patient's medication considering the circumstances of their state, it is necessary for diabetes patients. The study explores future research using Artificial Intelligence algorithms and presents a novel framework for monitoring the health of diabetes patients. However, technological problems and fluctuations may affect the accuracy of the information gathered by IoT devices, such wearables and sensors. The dependability of health information may be impacted by variables such as sensor reliability, measurement, and possible interference with the signal, which might have an effect on the monitoring system's general reliability [22].

In anticipation of the coronavirus disease-19 epidemic, audio-based telemedicine services for consultations and prescription medications were initially implemented in Korea. This study looked at how telehealth services affected health-care usage and drug prescription trends in hypertensive and diabetic individuals. The claims information from the Health Insurance Evaluation and Assessment Services for 2019 to 2021 were utilized. The difference-in-difference technique was utilized to compare the impact of telehealth treatments on the subjects as well as control groups before and following the period of intervention. Individuals in the untreated category employed in-person outpatient treatment, whereas those in the actual category received combined telemedicine and in-person treatments. The study comprised hypertensive individuals and diabetic patients. Telehealth services were linked to a rise in appointments with physicians among hypertensive. Patients with hypertension had a reduction in hospitals and visits to emergency rooms. In addition, policy execution has led to a rise in the medication possession ratio and the percentage of suitable prescriptions among patients with diabetes and high blood pressure. The data indicate a link between telemedicine service implementation and enhanced habits in health care usage and drug prescription, indicating telemedicine's potential utility in chronic illness management. However, the results of the investigation may be unique to the Korean health care system and not immediately relevant to other nations with differing healthcare facilities, laws, and cultural settings.

Factors such as regulations, consumer preferences, and healthcare professional practices might differ greatly between nations [23].

The literature review includes a range of subjects related to health monitoring systems, with a special emphasis on the treatment of diabetes patients and the use of technology to deliver healthcare more effectively. The first study uses Internet of Things devices to follow patient cases and offers an integrated voting classifier using neural networks for online diabetes monitoring. Although it shows better prediction rates, accuracy, recall, and F-measure than previous models, it also emphasizes possible limits by pointing out that data transfer to the cloud depends on continuous internet availability. The efficacy of E-Health treatments for glycaemic control in type 2 diabetes is investigated in the second research, which highlights the advantages of brief intervention durations particularly less than six months and the usefulness of SMS in lowering HbA1c levels. In discussing wearables' revolutionary potential in the healthcare industry, the third piece of literature focuses on how intelligently they can monitor a range of health indicators. The study does, however, recognize that demographic heterogeneity and device variety make it difficult to generalize outcomes. The fourth research addresses the healthcare requirements of the aging population by suggesting a multi-agent-based health monitoring system that provides enhanced services and mobility for chronic patients. Finally, the fifth research highlights the use of artificial intelligence algorithms for predictive modelling in the introduction of an Internet of Things-based health monitoring system for patients with diabetes. There are still issues with scalability and data dependability, despite encouraging developments. Together, these studies demonstrate the potential advantages of integrating technology into healthcare, but they also emphasize the necessity of resolving related issues in order to achieve widespread acceptance and effectiveness.

## III. PROBLEM STATEMENT

The present techniques used to monitor diabetic patients in electronic health systems frequently encounter issues such imprecise data analysis, restricted flexibility to accommodate different patient profiles, and inadequate real-time tracking capacity. By creating an advanced e-health System with an Integrated Intelligent Tracking Mechanism utilizing a Hybrid adaptive machine learning Classifier—more precisely, the suggested CNN-SVM model—this research seeks to address these problems. Accurate and individualized monitoring has been impeded by the present limitations in conventional approaches, which include their reliance on static thresholds and unsophisticated machine learning algorithms [6]. By combining the advantages of SVM for classification and CNN for feature extraction, the suggested CNN-SVM model outperforms other approaches in terms of robustness and efficiency when handling dynamic, complex patient information.

## IV. PROPOSED HYBRID ADAPTIVE MACHINE LEARNING CLASSIFIER

The methodology used in the research is an exhaustive approach to creating a sophisticated e-health system for

tracking diabetic patients, combining sophisticated surveillance techniques with an integrated adaptable machine learning classifier. The dataset includes information from sixty-seven consecutive days of exams for sixty-two diabetic patients, of which forty-four were male and eighteen were female. Using Min-Max normalization as a preprocessing step, the 13,173 concentrations of glucose information points and five attributes are scaled uniformly. PCA is used in feature extraction to reduce dimensionality and find important factors that influence patient features and fluctuations in glucose levels. The categorization and tracking mechanism, which uses a hybrid convolutional neural network integrated with a support vector machine, is the central component of the suggested system. Seven layers of CNN automatically identify important characteristics, while SVM guarantees robust categorization. The paper presents an innovative Hybrid CNN-SVM framework and demonstrates how SVM's competence in binary categorization and CNN's extraction of characteristics capabilities complement each other. Modelling parameter tuning is done using the Grey Wolf Optimization framework, which takes its information from the cooperative optimization procedure utilized by wolf packs. This optimization, inspired by nature, improves the model's efficiency. Fig. 1 depicts the general architecture of the suggested model.



Fig. 1. Overall structure of the proposed model.

### A. Data Collection

The sixty-two diabetic patients (forty-four men and eighteen women) whose medical histories required an average of sixty-seven days of testing were added to the database for this study. There are five characteristics and 13,173 glucose concentration information points in the glucose concentration sets [24].

### B. Preprocessing using Min-Max Normalization

Standardized scaling of characteristics was ensured by applying Min-Max normalization throughout the preparation of the dataset. 13,173 glucose concentration information points and five attributes are included in the collection of data. The level of glucose and other features were subjected to Min-Max normalization in order to scale the data within an acceptable range, usually [0, 1]. In order to enable more accurate and efficient modelling of diabetes-related variations and patterns within the information set, normalization is a crucial step in removing scale-related inefficiencies and verifying that each characteristic contributes proportionately to the study.

The research may scale the input information into an appropriate range by using Min-Max normalization, enabling a more rapid and precise evaluation. The relationships between the data sources are preserved without affecting the original information by making effort that all of the variables that are entered are transformed into a similar interval utilizing the normalization approach. The capacity of Min-Max normalization to maintain that connect between information points is one of its main advantages. This demonstrates that following normalization, the information's relative distribution and structure are still maintained. Preserving the fundamental patterns and correlations in the information is essential to accurately manage the elements influencing the projected values. Utilize Eq. (1) to illustrate the Min-Max normalizing method.

$$k = min_{new} + (max_{new} - min_{new}) * \left(\frac{k - min_r}{max_s - min_s}\right) \text{ (1)}$$

The starting measurement point is represented by "k" in this equation, the requisite values for the standardized data are denoted by $min_{new}$ and $max_{new}$, and the maximum and minimum values are indicated by $max_r$ and $min_r$, respectively. The speed and effectiveness of Min-Max normalization are its key benefits when handling an extensive amount of information points.

### C. Feature Extraction using Principal Component Analysis

An essential first step in evaluating high-dimensional datasets, like the diabetic patient monitoring dataset, is feature extraction. A common method for reducing dimensionality in information includes Principal Component Analysis, which tries to extract the most crucial information from the data being analyzed while removing less significant aspects. PCA can assist in determining the major factors that most influence changes in individual characteristics and glucose levels within the framework of diabetes surveillance.

The main elements, or eigenvectors of the initial information's covariance matrix, are the key idea that supports

principal component analysis. Let $Z$ represent the initial information matrix, which has $m$ characteristics and $m'$ observations. Eq. (2) calculates the covariance matrices $D$.

$$D = \frac{1}{m}(Z - \bar{Z})^T(Z - \bar{Z}) \qquad (2)$$

where, the mean-centered matrices is denoted by $\bar{Z}$.

Eq. (3) corresponds to D's eigenvalues λ and eigenvectors $w$.

$$Dw = \lambda w \qquad (3)$$

The primary elements are represented by these eigenvectors, and the quantity of variation they contain is shown by the associated eigenvalues. The primary elements are obtained in order of importance by categorizing the eigenvectors in decreasing order of eigenvalues.

Choose the highest $h$ eigenvectors that correlate to the $h$ greatest eigenvalues in order to minimize dimensionality. The most significant data in the information is captured by these $h$ main elements. Calculating the initial information $Z$ onto the chosen primary elements yields an updated characteristic matrix $X$, which is shown in Eq. (4).

$$X = ZW_h \qquad (4)$$

And the matrix containing the initial $h$ eigenvectors is denoted by $W_h$.

PCA makes dimensionality reduction simpler by maintaining the most crucial characteristics that have a substantial impact on the information's variance. As in the case of diabetes patients being monitored with various features and glucose concentration measurements, this reduction is especially useful in situations when the initial data set contains a large number of attributes. Because each principal component of the reduced-dimensional space produced by PCA maintains a mix of initial characteristics, interpretability is improved. Finding appropriate data for diabetes management is made easier by the more intuitively comprehension of patterns and linkages made possible by the information's decreased spatial visualization. PCA can identify the main factors that account for the majority of the variation in patient information when it comes to diabetes. For example, it might identify a variety of behavioural and physical characteristics that are essential to comprehending and tracking the course of the illness. More effective modelling and predictive evaluation may also benefit from the changed information in the reduced space. Comprehensive validation and efficacy assessment are essential to determine the effectiveness of PCA in the diabetic patient monitoring systems. This involves evaluating how well models constructed using the initial information performs in comparison to those developed using the PCA-transformed information, while taking consideration factors including computational effectiveness, interpretability of the models, and categorization accuracy. The evaluation's findings will provide information regarding the value and contribution of PCA to improving the efficiency of the system across all components.

## D. Classification and Tracking Mechanism using Hybrid CNN Networks with SVM

Hybrid CNN-SVM integration is used in the suggested Categorization and Tracking procedure, providing a comprehensive method of categorization and tracking. By utilizing two maximum-pooling layers, a pair of convolutional layers, and three layers that are interconnected for obtaining complex patterns from the input information, the combined approach builds on CNNs' advantages for automatic characteristic identification. The array of features is then loaded into an SVM to achieve reliable categorization. The predicted accuracy and flexibility of the model are improved by the collaboration. A thorough and precise categorization and tracking system is ensured by the CNN's automated identification of important characteristics and SVM's skill with high-dimensional spaces and intricate decision limits. The combination of these two effective algorithms operates in conjunction to effectively manage the dataset's complexities, creating a hybrid model that specializes at tasks including patient monitoring and illness prediction. Through exhaustive instruction, validation, and fine-tuning procedures, the efficacy of this hybrid technique is assessed, displaying its capacity as an innovative approach for difficult categorization problems.

*1) The CNN model:* The first classifier that this study suggests is a seven-layer CNN. The architecture consists of a pair of maximum-pooling layers, a pair of convolution layers, and three fully linked layers. Using a CNN instead of more traditional machine learning techniques has the main advantage that it can track and categorize important features without requiring human intervention. Its higher capacities need less human involvement due to its independence from pre-processing. Indeed, the different stages of convolution, that include a fixed number of filters, are used for autonomously obtaining the maps of attributes. $N'$ is the vectors' dimensions $r$, $n'$ and $k$ represent the vector's indicators, and $s = (s(n'))_{n'}$ is the result of the kernel's convolution of the input data $r = (r(k))_k$ As per Eq. (5), the stages of convolution integrate their vector inputs employing different filters. The maximum pooling layer reduces the complexity of networks by collecting the maximum values within a particular filter region. The complete classification is generated by the completely connected layer, which gathers data from the entire characteristic map. It frequently appears in the centre of the output layer.

$$r(n') = \sum_{k=0}^{N'-1} r(k) . f(n' - k) \qquad (5)$$

The initial layer receives as input a subsection with 500 points. This convolutional layer uses five 1×13 filters with a duration of 1 to convolve the various filters with the five-hundred-point values in accordance with Eq. (5). The result is five characteristic maps. The second layer is a maximal pooling layer with an initial pool size of 2 and a position offset of four. This layer reduces the size of the distinctive maps by combining a 1×2 filtering onto each of the previously created feature maps. As a result, the network must analyse less information and acquire fewer variables. Consequently, a

simulation is better able to manage variations in the location of input features. Applying a second layer of convolution with 10 filters that are each 1×9 in size and an advance of only one occurs next. This set of filters is used to extract higher-level features from dimensionally condensed maps of features. The fourth layer, which promotes pooling, carries out the characteristics and responsibilities of the initial layer. There are then three completely connected surfaces totalling 40, 20, and 2 properties. After being flattened, the output from the previous layer was utilized as the initial layer's input. With the exception for the highest layer, which employs the softmax activation function, other layers utilize the Leaky ReLU stimulating function.

Glorot uniform initialization is used for establishing the structure's weights, and backpropagation is employed for updating them over a maximum of sixty-four batches. The simulation is constructed throughout thirty epochs. Consider $\hat{q}_u$ as the estimated probability that the section u has diabetes, as found at the system's output. Eq. (6) shows how the binary cross-entropy function is used to quantify the simulation's loss in detecting the binary problem.

$$\mathcal{L}(q) = \frac{-1}{M} \sum_{u \in \mathfrak{I}_T} y_u . \log(\ \hat{q}_u) + (1 - y_u) . \log(1 - \ \hat{q}_u) \ (6)$$

When $M$ is the total number of these sections, $q$ is the epoch's indicators, and $\mathfrak{I}_T$ is the set of fragment indexes used for developing the system, then $M$ is the greatest number of $\mathfrak{I}_T$. The score that is calculated employing the cross-entropy represents the average deviation among the actual and projected values. The objective is to lower the score, where 0 represents the ideal cross-entropy.

The framework's range of variables can be modified through the use of grid-based searches and trial-and-error techniques. The tuning strategy grid search is used to find the optimal hyperparameter variables. It is a process that traverses over a manually selected subset of the targeted algorithm's hyperparameter space in detail. In this study, grid search is used to adjust the total amount of batches and epochs, and trial and error is used to determine the filter dimensions and durations.



Fig. 2. Structure of the suggested hybrid adaptive deep classifiers.

*2) The CNN-SVM model:* This section illustrates the recommended network's evolution. Instead of preserving the CNN network's last segment, which is responsible for classification, the study replaces it with an SVM classifier. Fig. 2 depicts the construction of the proposed hybrid CNN-SVM method. Algorithms that combine the outcomes of two or more distinct techniques are referred to as ensemble learning. Diabetes classification, monitoring, and early detection have all benefited from the use of collaborative learning in the field of healthcare. In short, a reduced CNN network is kept to extract attributes, and then the classification is done employing SVM. Consequently, a hybrid CNN-SVM technique is proposed for the diabetes tracking and classification datasets. The best features of CNN and SVM classifiers are combined in the proposed method. The trained CNN uses self-learning algorithms to identify the distinctive maps that are transmitted to the SVM for binary detection. CNN acts similarly to individuals and is particularly adept at remembering invariant local properties. It could extract the most unique information from the initial information. Support vector machines are classification techniques that learn to differentiate between input data's binary labels. Instances are used by a learning algorithm to educate it how to label objects. An SVM is simply a statistical approach that maximizes a particular statistical function with respect to a set of information. It finds a distinct hyperplane that divides information by extending a dataset's margins. The margin is the lowest length, split by a hyperplane, between two pieces of

data. The linear SVM technique has been extended to consider non-linear problems by projecting the data onto a higher-dimensional space. This method has proven to be quite successful because of how easily high-dimensional data can be handled and since linear strategies are simple to comprehend. The findings demonstrate that SVM responds well for binary tracking and categorization but inadequately for information with noise. Because of its basic architecture, SVM presents difficulties when learning deep properties. The hybrid CNN-SVM model proposed in this paper replaces the SoftMax layer of CNN with a non-linear SVM functioning as a binary classification algorithm.

### E. Grey Wolf Optimization Framework for Fine-tuning the Parameters

The novel method for adjusting parameters in this study is the Grey Wolf Optimization framework. GWO imitates the cooperative optimization process that occurs inside a wolf pack and is inspired by the social structure and hunting habits of grey wolves. By distributing the responsibilities of alpha, beta, and delta wolves in the framework of variable optimization, GWO dynamically modifies the stages of exploration and exploitation. While beta wolves investigate the areas within the findings made by alpha wolves, delta wolves concentrate on local research, while alpha wolves take the initiative in global exploration. The convergence towards the ideal values for parameters for the given position is facilitated by the collaborative and hierarchy optimization technique. The fundamental model's parameters are adjusted using the GWO structure, providing that the algorithm is sensitive to the unique features of the dataset as well as optimized for efficiency. By adding GWO, the fine-tuning procedure gains a sophisticated and nature-inspired component that increases the efficacy and efficiency of optimization of parameters for the intended application.

A suggested meta heuristic method is called GWO [25]. The method was influenced by the grey wolf killing strategy and pack structure. Grey wolves have a very hierarchical structure and socialize in packs. The leaders of the wolves, the alphas (α), now make all the decisions. Beta (β) wolves, which belong to the next level, help alpha wolves with their tasks. The final person, Omega (ω), is victimized in this system. A wolf is also known to as a delta (δ) wolf if it does not fall into any of the aforementioned classes. Grey wolves attempt to encompass a food source, assault, and kill, then explore for additional prey in accordance with this well-established structure. Wolves use hunting as a means of enclosing their prey, locating and killing animals, and engaging in conflict with their prey. Grey wolves on a hunting excursion circle their prey according to Eq. (7) and Eq. (8).

$$\vec{L} = |\vec{f} \cdot (\overrightarrow{U_z}(m) - \vec{U}(m)| \tag{7}$$

$$\vec{U}(m+1) = \overrightarrow{U_z}(m) - \vec{Q} \cdot \vec{L} \tag{8}$$

$\vec{Q}$ and $\vec{L}$ constitute efficient vectors with the subsequent definitions, which are presented in Eq. (9) and Eq. (10). Where $\vec{U}$ represents the location of the wolf in a circular

configuration, $\overrightarrow{U_z}$ is the vector position of the prey, and *mis* is the current time.

The wolf's position in a circular arrangement is represented by $\vec{U}$ in Eq. (9) and Eq. (10), while the prey's vectors position is represented by $\overrightarrow{U_z}$ in equations, and the present time is indicated by *m*, and the efficient vectors with matching definitions are $\vec{Q}$ and $\vec{L}$.

$$\vec{Q} = 2\vec{p} \cdot \vec{d_1} - \vec{p} \tag{9}$$

$$\vec{f} = 2 \cdot \vec{d_2} \tag{10}$$

The elements $\vec{d}$ and $\vec{d_2}$, where the component *d* is continuously decreasing from 2 to 0, contain random vectors evenly dispersed between 0 and 1. It has been suggested that the *α, β, and δ* wolves understand it easier since the exact spot of the food is never known in advance. Eq. (11), Eq. (12), and Eq. (13) are utilized to find the victim's location based on the locations of the wolves.

$$\vec{L}_\alpha = |\vec{f_1} \cdot \vec{U}_\alpha - \vec{U}|, \vec{L}_\beta = |\vec{f_2} \cdot \vec{U}_\beta - \vec{U}|, \vec{L}_\delta = |\vec{f_3} \cdot \vec{U}_\delta - \vec{U}| \tag{11}$$

$$\vec{U}_1 = \vec{U}_\alpha - \vec{Q}_1 \cdot \vec{U}_\alpha, \vec{U}_2 = \vec{U}_\beta - \vec{Q}_2 \cdot \vec{L}_\beta, \vec{U}_3 = \vec{U}_\delta - \vec{Q}_3 \cdot \vec{L}_\delta \tag{12}$$

$$\vec{U}(m+1) = \frac{\vec{U}_1 + \vec{U}_2 + \vec{U}_3}{3} \tag{13}$$

The next stage is to follow the victim (exploitation), if the study has an approximate position. The vector $\vec{Q}$ may be used to do this as the circumstance of wolves becomes nearer to the prey's location as *p* in Eq. (11) decreases from 2 to 0. Moreover, by removing the requirement for local averages, variables *f and Q* also contribute to maintaining the method's exploring capabilities. The accessibility of food and the difficulty of foraging may be altered by the variable *f*, but it can also affect a *Q* value larger than one, or *|Q| > 1*, which pushes the wolves to depart from their diet and search it out. Once the method is applied to a group of wolves for a set number of iterations, Eq. (13) will finally display the prey's position or the optimal region on Globe.

| Algorithm 1: SVM-CNN with GWO |
| --- |
| *// Input Data* |
| *// Assume patient data is a matrix where each row represents a patient's information* |
| *// Columns include glucose concentration data and other relevant attributes* |
| *patient_data = load_patient_data()* |
| *// Preprocessing* |
| *Normalized_data = min_max_normalization(patient_data)* |
| *// Feature Extraction using PCA* |
| *extracted_features = principal_component_analysis(normalized_data)* |
| *// Split data into training and testing sets* |
| *train_data, test_data = split_data(extracted_features)* |
| *// Build and Train CNN Model* |
| *cnn_model = build_and_train_cnn(train_data)* |
| *// Obtain CNN Output* |
| *cnn_output = get_cnn_output(cnn_model, extracted_features)* |
| *// Initialize SVM Parameters* |

```
svm_parameters = initialize_svm_parameters()
// Build and Train Hybrid CNN-SVM Model
hybrid_model = build_and_train_hybrid_model(cnn_output,
svm_parameters, train_data)

// Grey Wolf Optimization for Parameter Fine-tuning
optimized_parameters =
grey_wolf_optimization(hybrid_model.parameters, iterations=100)
// Build Final Model with Optimized Parameters
final_model = build_and_train_hybrid_model(cnn_output,
optimized_parameters, train_data)
// Evaluate Final Model
accuracy, precision, recall, f1_score = evaluate_model(final_model,
test_data)
```

## V. RESULTS AND DISCUSSION

The study uses an exhaustive approach to create an advanced e-health system that combines advanced surveillance techniques with an adaptable machine learning classifier to monitor diabetes patients. The dataset includes information from sixty-two diabetic patients who were examined over the course of sixty-seven days. The 13,173 glucose concentration measurement values and five characteristics are evenly scaled as a preprocessing step using Min-Max normalization. The next step is to extract characteristics using Principal Component Analysis, which reduces dimensionality and identifies important factors impacting patient characteristics and fluctuations in glucose levels. The main component of the suggested system is the tracking and classification mechanism, which makes use of a hybrid convolutional neural network coupled with a support vector machine. Significant characteristics are automatically identified by the seven-layer CNN, and robust classification is ensured by SVM. In order to highlight the complementing advantages of CNN's feature extraction skills and SVM's binary classification, the study presents novel hybrid CNN-SVM architecture. The Grey Wolf Optimization framework is modelled after the cooperative optimization process seen in wolf packs, is used to tune the model's parameters and improve its efficiency.

### A. Performance Evaluation

Evaluation indicators are essential for assessing the effectiveness of classification. A calculation of accuracy is the approach that is most commonly utilized for this goal. How well a classifier identifies sample datasets may be used to measure how accurate it is for any given collection of information. Because depending entirely on the accuracy measure will prevent you from making the best assessments conceivable. The researchers also used other parameters to assess the classifier's effectiveness. Measures of accuracy, recall, precision, and F1-score were employed to assess the effectiveness of the proposed method. The definitions of each metric are described as follows

- $T_{pos}$ (True Positive) refers to the amount of information that has been correctly categorized.

- The term $Fpos$ (False Positive) represents the volume of reliable information that was incorrectly categorized.

- False negatives ($F_{neg}$) are instances where incorrect information has been given an actual classification.

- The categorization of incorrect information values is referred to as $T_{neg}$ (True Negative).

*1) Accuracy:* The accuracy of the classifier shows how often it makes the correct prediction. Accuracy is defined as the ratio of correct estimations to all other reasonable theories. It is demonstrated by Eq. (14).

$$Accuracy = \frac{T_{pos}+T_{neg}}{T_{pos}+T_{neg}+F_{pos}+F_{neg}} \quad (14)$$

*2) Precision:* Evaluating a classifier's precision, or degree of accuracy, yields the number of possibilities that are properly identified. Increased reliability leads to fewer false positives, but lower precision results in many more. Precision is defined as the proportion of properly classified cases relative to all occurrences. It is defined by Eq. (15).

$$P = \frac{Tpos}{Tpos+Fpos} \quad (15)$$

*3) Recall:* Recall determines a categorization's sensitivity, or how much pertinent information it generates. As recollection improves, $F_{neg}$ total amount decreases. The percentage of occurrences that have been accurately classified to all of the expected instances is called recall. This is demonstrable by Eq. (16).

$$R = \frac{Tpos}{Tpos+Fneg} \quad (16)$$

*4) F1-Score:* Addition of precision and recall yields an association of measurements known as the F-measure, which represents the weighted average of accuracy and recall. It is characterised by Eq. (17).

$$F1 - score = \frac{2 \times precision \times recall}{precision \times recall} \quad (17)$$

*5) ROC Curve:* In deep learning and machine learning, area under the ROC curve, or AUC, is a popular assessment statistic for binary categorization issues. The area under the curve (AOC) is a visual depiction of the receiver operating characteristic (ROC) curve that shows how effective the binary identification technique is. In a binary categorized issue, the classifier determines whether the incoming data is part of a positive or negative partition. The ROC curve displays the $Tpos$ vs. the $F_{pos}$ for different categorization parameters. AOC values range from 0 to 1, with higher numbers denoting more efficiency. An optimum classifier has an AOC of one, whereas a totally randomly assigned classifier has an AOC of 0.5. Since the approach takes into account every conceivable level of detection and offers only one statistic for comparing the effectiveness of various classifiers.

Fig. 3.    Training and testing accuracy.

The training and testing accuracy scores at different epochs of the model training procedure are shown in Fig. 3. Training and testing accuracy exhibit a steady rising trend with an increase in training epochs, suggesting that the model is performing better. The model obtains a testing accuracy of 75% and a training accuracy of 76.8% at the beginning of training.



Fig. 4.    Training and testing loss.

The training accuracy increases with the number of epochs, reaching at 99% after 90 and 99.6% after 100 epochs. A similar pattern may be found in the associated testing accuracy, which shows how well the model generalizes to new information. The model's ability to learn from and adapt to the dataset is seen by the significant rise in accuracy from 10 to 100 epochs. The final epochs achieve a high degree of accuracy, indicating a resilient and well-trained model. The model is more reliable in correctly predicting patterns associated to diabetes because training and testing accuracy converge at higher epochs, indicating efficient learning without overfitting. The training and testing loss values at various epochs during the algorithm's training procedure are shown in Fig. 4. Lower numbers indicate higher model performance. Loss values show the difference between the expected and actual values. Training and testing loss are both rather high in the early phases of training (at 10 epochs),

indicating the model's inadequate ability for precise result prediction. On the other hand, training and testing loss consistently decrease with the number of epochs, indicating enhanced model convergence and accuracy in predictions. The training loss dramatically drops to 0.06 by the 100th epoch, demonstrating that the system effectively eliminates mistakes throughout the learning process from the training set. Additionally, the testing loss drops to 0.14, indicating that the model can generalize even on untested information. The model's capacity for identifying diabetes-related characteristics is supported by the consistent reduction in loss values over epochs, which shows effective learning, and the similarity of testing and training losses, which suggests the model, maintains high accuracy without overfitting.

TABLE I.        PROPER AND IMPROPER CATEGORIZED INFORMATION

| Methods | SVM [26] | DT [26] | SMO [27] | CNN-SVM |
|---|---|---|---|---|
| Improper Categorized Data | 8.482% | 11.030% | 13.715% | 0.312% |
| Proper Categorized Data | 89.115% | 84.541% | 79.455% | 99.851% |

With a focus on data pertaining to diabetes, Table I and Fig. 5 examine how well various categorization techniques performed in terms of accurate and incorrect information categorization. Support Vector Machine (SVM), Decision Trees (DT), Sequential Minimal Optimization (SMO), and the suggested hybrid approach, CNN-SVM, are among the techniques assessed. The rates at which each approach misclassifies information are indicated by the percentages in the Improper Categorized Information. The percentages of incorrectly categorized information for SVM, DT, and SMO are greater (8.482%, 11.030%, and 13.715%, respectively). The CNN-SVM hybrid strategy, on the other hand, performs noticeably better than these techniques, attaining a very low rate of 0.312% in incorrect classification.



Fig. 5.    Proper and improper categorized information.

The Proper Categorized Information, on the other hand, demonstrates how accurately each approach categorizes information pertaining to diabetes. This is where CNN-SVM excels, exceeding SVM, DT, and SMO, with percentages of 89.115%, 84.541%, and 79.455%, respectively, in appropriate categorization with an astounding 99.851% accuracy. The outcomes demonstrate the suggested CNN-SVM hybrid model's improved performance in classifying diabetes-related information with accuracy, indicating its potential as a useful method for patient monitoring and illness prediction.



Fig. 6.    Fitness improvement over iterations.

The Grey Wolf Optimization algorithm's enhancement in efficiency as iteratively refines the model's parameters is demonstrated in Fig. 6 by the Fitness Improvement over Iterations. The fitness value which indicates the optimization of the objective function increases during the first iterations while the algorithm searches the parameter space. A discernible increase in fitness is shown as the iterations continue on, suggesting that the GWO method is effective in fine-tuning the parameters to get improved placement with the optimization objective. Fitness values show a constant decreasing trend, which indicates that the algorithm is efficient in converging towards the best parameter combinations. The above chart provides a visual demonstration of the GWO algorithm's capacity to dynamically modify parameters to improve the model's efficiency continually. It also demonstrates the algorithm's effectiveness in fine-tuning for optimal outcomes across a number of rounds.



Fig. 7.    ROC curve.

Fig. 7 shows a plot of True Positive Rate (Sensitivity) vs. False Positive Rate (1 - Specificity) over several threshold values for a binary categorization model, which represents the Receiver Operating Characteristic (ROC) Curve. The True Positive Rate progressively improves as the discriminating threshold reduces from 0.6 to 0.6 in the given figure of threshold values and associated False Positive Rates, indicating the model's capacity to accurately detect positive events. Additionally, there is an increase in the False Positive Rate, which signifies the occurrences of the model misclassifying negative cases as positive. The relationship between sensitivity and specificity is graphically represented by the ROC Curve, which offers an understanding of the model's overall discriminating strength over a variety of threshold values.

TABLE II.        COMPARISON OF ROC OF PROPOSED METHOD WITH OTHER EXISTING APPROACHES

| Models | Levels |
|---|---|
| SVM | 0.8148 |
| DT | 0.8052 |
| SMO | 0.8264 |
| Proposed CNN-SVM | 0.8687 |

Table II presents a comparison of the Receiver Operating Characteristic (ROC) performance metrics for different models, specifically Support Vector Machine (SVM), Decision Tree (DT), Sequential Minimal Optimization (SMO), and the proposed method, a Convolutional Neural Network-Support Vector Machine hybrid (CNN-SVM). The ROC values serve as indicators of the models' ability to discriminate between classes, with higher values suggesting better performance. In this context, the proposed CNN-SVM demonstrates the highest ROC value of 0.8687, indicating superior discriminative capabilities compared to SVM (0.8148), DT (0.8052), and SMO (0.8264). The results suggest that the hybrid approach, combining Convolutional Neural Network and Support Vector Machine, outperforms traditional machine learning models in the specific task or dataset under consideration, emphasizing its potential for enhanced predictive accuracy and classification performance.

The suggested CNN-SVM model's performance metrics for diabetes prediction are compiled in Fig. 8, which displays outstanding outcomes for all major assessment parameters. With an impressive 99.6% accuracy rate, the model demonstrates its ability to accurately categorize occurrences. With a remarkable 99.4% precision rate a measure of the model's accuracy in positive predictions it is clear that there is little chance of false positives. The model's strong sensitivity is further demonstrated by the recall metric, which measures the model's capacity to identify all positive events and is now 99.4%. At a remarkable 99.5%, the F1-Score a balanced metric of accuracy and recall highlights the CNN-SVM model's overall efficacy in diabetes prediction. All of these findings indicate the suggested model's stability and dependability, highlighting its possibilities as an innovative technology for precise and effective diabetic patient monitoring.

Fig. 8.   Model performance.

TABLE III.   COMPARISON OF PERFORMANCE METRICS OF PROPOSED METHOD WITH OTHER EXISTING APPROACHES

| Models | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| SVM | 96.45 | 95.11 | 95.15 | 95.17 |
| DT | 95.34 | 95.23 | 95.23 | 95.10 |
| SMO | 97.12 | 96.09 | 96.10 | 96.10 |
| Proposed CNN-SVM | 99.6 | 99.4 | 99.4 | 99.5 |

The suggested CNN-SVM model and other methods, such as Support Vector Machine (SVM), Decision Tree (DT), and Sequential Minimal Optimization (SMO), are extensively contrasted using performance metrics in Table II and Fig. 9. With an outstanding 99.6% accuracy rate, the suggested CNN-SVM model performs better than alternative approaches. The suggested model has outstanding performance as seen by its 99.4%, 99.4%, and 99.5% precision, recall, and F1-Score metrics. By contrast, SVM attains competitive precision, recall, and F1-Score values in addition to a high accuracy of 96.45%.

Accuracy ratings of 95.34% and 97.12% for DT and SMO, respectively, also show satisfactory results. However, the suggested CNN-SVM model performs better overall on all assessed parameters, highlighting its effectiveness in tasks involving the prediction and categorization of diabetes. These findings demonstrate the hybrid CNN-SVM approach's potential for innovative healthcare applications, especially when it comes to diabetic patient monitoring.

*B. Discussion*

By combining a variety of sophisticated surveillance approaches with an adaptive machine learning classifier, the study described here provides a comprehensive method for creating an advanced e-health system for diabetic patient monitoring. The research makes use of a dataset that includes sixty-two diabetic patients who were followed up on for sixty-seven days in a row. The dataset undergone significant preprocessing, which included feature extraction using Principal Component Analysis (PCA) and Min-Max

normalization. The fundamental component of the suggested system combines a support vector machine (SVM) with a hybrid convolutional neural network (CNN) for tracking and categorization. Grey Wolf Optimization framework is used to tune model parameters. Accuracy, precision, recall, F1-Score, and the Receiver Operating Characteristic (ROC) Curve are all used in the model's performance evaluation to give a thorough assessment of its prediction competencies. The results show that the model is adaptable in classifying data relevant to diabetes and show a notable improvement in accuracy throughout training epochs, reaching an astonishing 99.6%. The proposed CNN-SVM model has potential for accurate and efficient diabetic patient monitoring by outperforming other conventional methods like SVM [26], DT [26], and SMO [27]. The Fitness Improvement over Iterations graph, which illustrates the study's findings, provides an understanding of how the Grey Wolf Optimization method affects the model's effectiveness. This graph shows how the method refines parameters iteratively, improving fitness values and optimizing the objective function. The model's ability to learn and generalize well is further demonstrated by the continuous decline in loss values during training epochs. The suggested CNN-SVM model is superior to other current techniques in terms of accuracy, precision, recall, and F1-Score, as demonstrated by the comparison shown in Table II and Fig. 9. All of these results point to the possibility of creative healthcare applications using the hybrid CNN-SVM architecture and Grey Wolf Optimization algorithm, especially for accurate diabetes patient prediction and monitoring. The work demonstrates the possibility of combining machine learning and optimization approaches for better healthcare outcomes in addition to adding to the knowledge of diabetes patient monitoring.



Fig. 9.   Comparison of performance metrics of proposed method with other existing approaches.

## VI.   CONCLUSION AND FUTURE WORKS

In conclusion, this study has showcased an innovative method for monitoring diabetic patients, which has resulted in the creation of a sophisticated computerized health system that combines a sophisticated tracking system with a hybrid adaptive machine learning classifier. The system, which is

trained and optimized utilizing the Grey Wolf Optimization technique, takes advantage of the synergies between support vector machines (SVM) and convolutional neural networks (CNN) in hybrid architecture. The comprehensive assessment of conventional performance measures has proven the enhanced accuracy, precision, recall, and F1-Score of the suggested CNN-SVM model, indicating its efficacy in classifying data pertaining to diabetes. Analyses that compare the hybrid model to more conventional techniques like SVM, Decision Trees, and Sequential Minimal Optimization highlight the significant improvement in accuracy that the hybrid model offers. In addition, the study has provided informative visuals which provide an extensive comprehension of the learning dynamics and optimization of the model. These visualizations include fitness increase over iterations, ROC curves, training and testing accuracy graphs, loss curves, and more. The suggested system's resilience and ability to provide accurate and efficient diabetes patient monitoring emphasize its importance in improving e-health applications and creating opportunities for customized and adaptable healthcare solutions. The research makes a significant contribution by presenting a novel framework that combines machine learning, intelligent tracking, and optimization techniques. This framework paves up the opportunity for novel approaches to diabetes care in the e-health era. The generalizability and practicality of the model will be improved in subsequent work by extending the dataset to incorporate more varied demographic information and taking real-world deployment issues into account. Further research into the incorporation of real-time feedback from patient's mechanisms and the possibility of using edge computing to lower monitoring process latency might enhance the responsiveness and user involvement of the suggested e-health system.

## REFERENCES

[1] J. A. Andersen, D. Scoggins, T. Michaud, N. Wan, M. Wen, and D. Su, "Racial disparities in diabetes management outcomes: evidence from a remote patient monitoring program for type 2 diabetic patients," Telemed. e-Health, vol. 27, no. 1, pp. 55–61, 2021.

[2] H.-J. Seo, S. Y. Kim, S.-S. Sheen, and Y. Cha, "e-Health Interventions for Community-Dwelling Type 2 Diabetes: A Scoping Review," Telemed. E-Health, vol. 27, no. 3, pp. 276–285, 2021.

[3] L. Nachabe, R. Raiyee, O. Falou, M. Girod-Genet, and B. ElHassan, "Diabetes Mobile Application as a Part of Semantic Multi-Agent System for e-health," in 2020 IEEE 5th Middle East and Africa Conference on Biomedical Engineering (MECBME), IEEE, 2020, pp. 1–5.

[4] T. L. Michaud et al., "Association between weight loss and glycemic outcomes: A post hoc analysis of a remote patient monitoring program for diabetes management," Telemed. e-health, vol. 26, no. 5, pp. 621–628, 2020.

[5] S. Mangal et al., "e-Health initiatives for screening and management of diabetes in rural Rajasthan," Int. J. Diabetes Dev. Ctries., pp. 1–6, 2022.

[6] D. L. G. Rodrigues, G. S. Belber, I. da C. Borysow, M. A. Maeyama, and A. P. N. M. de Pinho, "Description of e-Health initiatives to reduce chronic non-communicable disease burden on Brazilian health system," Int. J. Environ. Res. Public. Health, vol. 18, no. 19, p. 10218, 2021.

[7] I. Smokovski and I. Smokovski, "Benefits of Centralized e-Health System in Diabetes Care," Manag. Diabetes Low Income Ctries. Provid. Sustain. Diabetes Care Ltd. Resour., pp. 73–83, 2021.

[8] R. J. Middelbeek, M. F. Bouchonville, S. Agarwal, and G. R. Romeo, "Application of telehealth to diabetes care delivery and medical training: challenges and opportunities," Front. Endocrinol., vol. 14, p. 1229706, 2023.

[9] F. Alanazi, V. Gay, R. Alturki, and others, "Modelling Health Process and System Requirements Engineering for Better e-health Services in Saudi Arabia," Int. J. Adv. Comput. Sci. Appl., 2021.

[10] S. Salehi, A. Olyaeemanesh, M. Mobinizadeh, E. Nasli-Esfahani, and H. Riazi, "Assessment of remote patient monitoring (RPM) systems for patients with type 2 diabetes: a systematic review and meta-analysis," J. Diabetes Metab. Disord., vol. 19, pp. 115–127, 2020.

[11] W. Zhang, B. Cheng, W. Zhu, X. Huang, and C. Shen, "Effect of telemedicine on quality of care in patients with coexisting hypertension and diabetes: a systematic review and meta-analysis," Telemed. e-Health, vol. 27, no. 6, pp. 603–614, 2021.

[12] A. K. Srivastava, Y. Kumar, and P. K. Singh, "A rule-based monitoring system for accurate prediction of diabetes: monitoring system for diabetes," Int. J. e-health Med. Commun. IJEHMC, vol. 11, no. 3, pp. 32–53, 2020.

[13] J. S. Coombes et al., "Personal activity intelligence (PAI) e-health program in people with type 2 diabetes: a pilot randomized controlled trial," Med Sci Sports Exerc, 2021.

[14] S. J. Fonda, S.-E. Bursell, D. G. Lewis, D. Clary, D. Shahon, and M. B. Horton, "The Indian health service primary care-based teleophthalmology program for diabetic eye disease surveillance and management," Telemed. e-health, vol. 26, no. 12, pp. 1466–1474, 2020.

[15] T. L. Michaud, J. Ern, D. Scoggins, and D. Su, "Assessing the impact of telemonitoring-facilitated lifestyle modifications on diabetes outcomes: a systematic review and meta-analysis," Telemed. e-health, vol. 27, no. 2, pp. 124–136, 2021.

[16] R. Xu, M. Xing, K. Javaherian, R. Peters, W. Ross, and C. Bernal-Mizrachi, "Improving HbA1c with glucose self-monitoring in diabetic patients with EpxDiabetes, a phone call and text message-based telemedicine platform: a randomized controlled trial," Telemed. e-health, vol. 26, no. 6, pp. 784–793, 2020.

[17] E. Kirkland et al., "Patient demographics and clinic type are associated with patient engagement within a remote monitoring program," Telemed. e-health, vol. 27, no. 8, pp. 843–850, 2021.

[18] P. EP et al., "Implementation of Artificial Neural Network to Predict Diabetes with High-Quality Health System," Comput. Intell. Neurosci., vol. 2022, 2022.

[19] X. Zhang et al., "Effects of e-health-based interventions on glycemic control for patients with type 2 diabetes: a Bayesian network meta-analysis," Front. Endocrinol., vol. 14, p. 1068254, 2023.

[20] T. Sivani and S. Mishra, "Wearable Devices: Evolution and Usage in Remote Patient Monitoring System," in Connected e-health: Integrated IoT and Cloud Computing, Springer, 2022, pp. 311–332.

[21] M. Humayun, N. Z. Jhanjhi, A. Almotilag, and M. F. Almufareh, "Agent-based medical health monitoring system," Sensors, vol. 22, no. 8, p. 2820, 2022.

[22] S. K. Polu, "Design of IoT Based Health Monitoring System for Diabetic Patients".

[23] M. Cho et al., "The Effect of Telehealth on Patterns of Health Care Utilization and Medication Prescription in Patients with Diabetes or Hypertension During COVID-19: A Nationwide Study," Telemed. e-health, 2024.

[24] A. Rghioui, A. Naja, J. L. Mauri, and A. Oumnad, "An IoT based diabetic patient monitoring system using machine learning and node MCU," in Journal of Physics: Conference Series, IOP Publishing, 2021, p. 012035.

[25] S. Mirjalili, S. M. Mirjalili, and A. Lewis, "Grey Wolf Optimizer," Adv. Eng. Softw., vol. 69, pp. 46–61, Mar. 2014, doi: 10.1016/j.advengsoft.2013.12.007.

[26] B. Godi, S. Viswanadham, A. S. Muttipati, O. P. Samantray, and S. R. Gadiraju, "e-healthcare monitoring system using IoT with machine learning approaches," in 2020 international conference on computer science, engineering and applications (ICCSEA), IEEE, 2020, pp. 1–5.

[27] A. Rghioui, J. Lloret, S. Sendra, and A. Oumnad, "A smart architecture for diabetic patient monitoring using machine learning algorithms," in Healthcare, MDPI, 2020, p. 348.

# Feature Selection Model Development on Near-Infrared Spectroscopy Data

## Case Study of Beef Freshness Quality Prediction

Ridwan Raafi'udin[1], Y. Aris Purwanto[2], Imas Sukaesih Sitanggang[3], Dewi Apri Astuti[4]

Department of Computer Science, IPB University, Bogor, Indonesia[1, 3]
Mechanical and Biosystem Engineering Department, IPB University, Bogor, Indonesia[2]
Department of Nutrition and Feed Technology, IPB University, Bogor, Indonesia[4]
Department of Informatics, Universitas Pembangunan Nasional Veteran Jakarta, Indonesia[1]

*Abstract*—**This study aims to develop a feature selection model on Near-Infrared Spectroscopy (NIRS) data. The object used is beef with six quality parameters: color, drip loss, pH, storage time, Total Plate Colony (TPC), and water moisture. The prediction model is a Random Forest Regressor (RFR) with default parameters. The feature selection model is carried out by mapping spectroscopic data into line form. The collection of lines is made into one line by finding the mean value. Next, apply the line simplification method based on angle elimination, starting from the smallest angle to the largest. Each iteration will eliminate one corner, reducing one column of data in the corresponding dataset. Then, the predicted value in the form of R2 will be collected, and the highest value will be considered the best feature selection formation. RFR prediction results with R2 values are as follows: color R2= 0.597, drip loss R2=0.891, pH R2=0.797, storage time R2=0.889, TPC R2=0.721, and water moisture R2=0.540. Meanwhile, after applying the feature selection model, the R2 values for all parameters increased to color R2=0.877, drip loss R2=0.943, pH R2=0.904, storage time R2=0.917, TPC R2=0.951, and water moisture R2=0.893. Based on the results of increasing the R2 value of the six parameters, an average value of increasing prediction accuracy of 17.49% can be taken. So, the feature selection method based on line simplification with an angle elimination system can provide very good results.**

*Keywords—Beef quality prediction; feature selection; machine learning; Random Forest Regressor*

## I. INTRODUCTION

People have consumed a lot of meat in the last few decades, and consumption has increased in recent years [1]. Beef is an alternative commodity that is widely consumed to meet the need for protein in many countries [2]. However, meat food products are products that rot quickly, especially under certain conditions, which can accelerate microbial growth [3]. Many cases of consumers getting sick are caused by microbes found in beef that are large in number or above standard [4].

The condition of the meat can change quickly, so a method is needed that can determine the current state of beef quality. One fast method for determining meat quality is to use near-infrared (NIR) technology [5]. NIR spectrometers can predict meat quality parameters, including chemical parameters, technological parameters, quality traits, fatty acids, and many mineral contents [6][7][8].

To speed up and simplify the process of determining meat quality, a portable or handheld NIR device can be built that can be taken anywhere [9][10]. The development of a portable NIR device can be applied in industries that require the process of determining the characteristics of meat products [11][12].

Machine learning can be used as a model to predict meat quality parameters [13]. Machine learning is also able to predict meat quality parameters, including color, tenderness, juiciness, and flavor [14]. The random forest (RF) algorithm works well and produces high accuracy in classifying cattle breeds [15]. Random Forest Regressor (RFR) performs well in predicting pH values in beef in real-time in a beef freshness monitoring system [16].

In this study, we propose a feature selection model on spectroscopic data to increase the accuracy of meat quality predictions. The beef quality parameters in question are color, drip loss, pH, storage time, total plate colony (TCP), and water moisture. The algorithm that will be used is RFR, with algorithm parameters by default. Based on the experimental results, the proposed method is able to increase accuracy in predicting the freshness quality of beef compared to results without using feature selection.

## II. LITERATURE SURVEY

### A. Feature Selection

Feature selection involves exploring algorithms designed to decrease the data's dimensionality, thereby enhancing the performance of machine learning. In datasets with N features and M dimensions, the primary goal is to minimize M to $M'$, where $M'$ is less than or equal to M [17]. Typically, feature selection often results in improved learning outcomes, such as increased accuracy in learning, reduced computational expenses, and enhanced interpretability of the model.

Recently, experts in fields like computer vision and text mining have introduced numerous feature selection methods. Through both theoretical frameworks and practical experiments, they've demonstrated the effectiveness of their approaches [18]. Feature selection using established line simplification methods such as the Ramer–Douglas–Peucker algorithm and Visvalingam on NIRS data encountered

problems in determining epsilon values, making it difficult to produce optimal datasets.

This study will compare the results of prediction accuracy from the proposed feature selection model with the established method from Scikit-Learn, namely SelectFromModel [19]. Comparison of results in the form of accuracy of value prediction on six meat quality parameters using R-squared values (R2) as part of model evaluation both for evaluation of the SelectFromModel feature selection model and for the proposed feature selection model.

*B. Random Forest Regressor*

Random forests consist of a collection of tree predictors where each tree relies on a randomly sampled vector, independent and identically distributed across all trees within the forest. As the number of trees in the forest increases, the generalization error for forests gradually converges to a limit. This error is influenced by both the effectiveness of the individual trees in the forest and the level of correlation among them [20].

The advantages of feature selection and its relevance to enhancing the effectiveness and interpretability of machine learning algorithms are well-established. In this context, we focus on incorporating feature selection into a Random Forest setup. We propose a method that integrates hypothesis testing with an estimation of the anticipated impact of an irrelevant feature while constructing a Random Forest [21].

The utilization of $R^2$ in this research for model assessment is due to its superior informativeness and reliability compared to SMAPE, along with its absence of interpretational constraints seen in metrics like MSE, RMSE, MAE, and MAPE[22]. The R2 value range is easy to understand because the value ranges between 0 and 1 and is indicated by a decimal number. The number 0 indicates poor model performance, and one or close to it indicates good model performance. The formula for calculating the R2 value can be seen in Eq. (1) [22].

$$R^2 = 1 - \frac{sum\ squared\ regression\ (SSR)}{total\ sum\ of\ squares\ (SST)}$$

$$= 1 - \frac{\sum(y_1 - \hat{y}_1)^2}{\sum(y_1 - \bar{y})^2} \qquad (1)$$

*C. Near-Infrared Spectroscopy*

Near-infrared spectroscopy (NIRS) holds a distinct position within the realm of bioscience and its associated fields, differing in characteristics and potential applications from infrared (IR) or Raman spectroscopy. This type of vibrational spectroscopy uncovers molecular details within a sample by detecting absorption bands arising from overtones and combined excitations [23].

The capacity of near-infrared reflectance spectroscopy (NIRS) to differentiate between Normal and DFD (dark, firm, and dry) beef and forecast quality characteristics within 129 Longissimus thoracic (LT) samples derived from three Spanish pure breeds: Asturiana de los Valles (AV; n = 50), Rubia Gallega (RG; n = 37), and Retinta (RE; n = 42). The outcomes obtained using partial least squares-discriminant analysis (PLS-

DA) demonstrated successful differentiation between Normal and DFD meat samples from AV and RG (with sensitivity exceeding 93% for both and specificity of 100% and 72%, respectively), whereas results for RE and the overall sample sets were less accurate. Soft independent modeling of class analogies (SIMCA) revealed 100% sensitivity for DFD meat across all sample sets (total, AV, RG, and RE) and over 90% specificity for AV, RG, and RE, but notably lower for the total sample set (19.8%). Utilizing NIRS quantitative models via partial least squares regression (PLSR) allowed dependable prediction of color parameters (CIE L*, a*, b*, hue, chroma). The findings from both qualitative and quantitative analyses hold promise for early decision-making in the meat production process to prevent financial losses and food wastage [24].

Growing apprehensions regarding contaminated meat have spurred the industry to explore novel, non-invasive techniques for swift and precise assessment of meat quality. The primary chromophores in meat (such as myoglobin, oxy-myoglobin, fat, water, and collagen) exhibit similar absorption patterns within the visible to near-infrared (NIR) spectral range. Consequently, variations in the structure and composition of meat can result in proportional disparities in light absorption [25].

### III. MATERIAL

*A. Sample Preparation*

This study used fresh beef objects obtained from traditional markets in Bogor City, West Java, Indonesia, and then brought to the laboratory using an ice box, as shown in Fig. 1. The part of the carcass used is tenderloin. The main sample used weighs 1 kg, as shown in Fig. 2.

The next step is to cut it into eight pieces of large samples weighing 17+2 grams and eight pieces of small samples of 3+1 grams, as shown in Fig. 3. All samples were placed on the laboratory table in Petri dishes and supported by wire gauze, as shown in Fig. 4. Large one of samples were used for testing the parameters of color, pH, WHC, water content, and NIR, while small one of samples were used for TPC measurements. The total samples used were eight samples per day and repeated for ten days to obtain 80 samples.



Fig. 1. Sample in the icebox.



Fig. 2. 1 Kg tenderloin.

Fig. 3.   Small pieces of sample.


Fig. 4.   The samples are placed on a petri dish with a wire mesh base.

## B. Data Acquisition

This study uses six beef freshness parameters, where each parameter uses a different tool. The six freshness parameters include color, drip loss, pH, storage time, TPC, and water moisture. The tool used to retrieve color data is Chromameter, as shown in Fig. 5.

The drip loss parameter value is obtained by measuring the weight of the sample between the initial and final times. The measurement interval is one hour with a total time span of seven hours so that eight drip loss data are produced starting from the 0th hour to the 7th hour. The tool for measuring sample weight is a laboratory scale, as shown in Fig. 6.


Fig. 5.   Measurement results using a chromameter.


Fig. 6.   Weight measurement process.

TPC parameter values were obtained from other laboratories and measured professionally by a third party. The samples sent are the same pieces used for measuring other parameters. Examples of samples are shown in Fig. 7.


Fig. 7.   Sample for TPC measurement.


Fig. 8.   pH measurement process.

The pH value parameter of the sample was measured using an electronic pH meter, as shown in Fig. 8. The storage time parameter value is obtained by simply storing the sample on the table from the 0th hour to the 7th hour with a break every hour. Then, the water moisture value is carried out using the drying or thermogravimetric method in an oven at a temperature of 105 degrees Celsius for 16 hours. Laboratory testing data collection activities were carried out for ten days, resulting in a dataset for 80 data rows, as shown in Table I.

TABLE I.        DATA EXAMPLE FROM LABORATORY

(A)

| Code | | | WHC | | | |
|---|---|---|---|---|---|---|
| *Day* | *Hour* | *Code* | *W0* | *Wt* | *ΔW* | *%drip loss* |
| 4 | 0 | d4h0 | 16.67 | 16.67 | 0.00 | 0% |
| 4 | 1 | d4h1 | 18.13 | 17.48 | 0.65 | 4% |
| 4 | 2 | d4h2 | 18.02 | 16.93 | 1.09 | 6% |
| 4 | 3 | d4h3 | 16.01 | 14.46 | 1.55 | 10% |
| 4 | 4 | d4h4 | 18.06 | 16.06 | 2.00 | 11% |
| 4 | 5 | d4h5 | 17.52 | 15.19 | 2.33 | 13% |
| 4 | 6 | d4h6 | 18.41 | 15.34 | 3.07 | 17% |
| 4 | 7 | d4h7 | 17.52 | 14.37 | 3.15 | 18% |

(B)

| Color | | | pH | Water Moisture (WM) | | | |
|---|---|---|---|---|---|---|---|
| *L* | *a* | *b* | | *Wdish* | *Wsample* | *Wt* | *%WM* |
| 33,05 | 16,23 | -1,42 | 5,51 | 25,47 | 3,45 | 26,32 | 75,36% |
| 33,92 | 17,51 | 1,28 | 5,46 | 29,60 | 3,93 | 30,50 | 77,10% |
| 36,46 | 16,79 | 0,08 | 5,55 | 28,46 | 3,84 | 29,41 | 75,26% |
| 33,29 | 16,04 | 3,25 | 5,50 | 28,81 | 4,00 | 29,82 | 74,75% |
| 32,47 | 15,68 | 4,94 | 5,55 | 30,16 | 3,92 | 31,19 | 73,72% |
| 29,50 | 11,18 | 3,00 | 5,50 | 29,13 | 3,94 | 30,19 | 73,10% |
| 31,62 | 12,67 | 3,23 | 5,48 | 25,12 | 3,67 | 26,16 | 71,66% |
| 28,29 | 10,88 | 5,43 | 5,59 | 27,97 | 3,68 | 28,97 | 72,83% |

$W_0$. initial weight; Wt. new weight; $\Delta W$. weight changes; L; lightness; a. red/green value; b. yellow/blue value; Wdish, the weight of the empty dish; Wsample, the weight of the current sample; Wt, weight changes.

## C. NIRS Data and Dataset

NIRS data is obtained by scanning samples using a NIRS sensor. The sensor used is the NeoSpectra Development Kit [26] as shown in Fig. 9 and Fig. 10. For the wireless data acquisition process, use a notebook unit equipped with requesting software that is available from the sensor manufacturer as shown in Fig. 11.



Fig. 11. Acquisition using NeoSpectraKit in progress covered by a box.

Measurement results in a spreadsheet file. The example of the averaged data can be seen in Table II. The original data of NIR Spectroscopy plotted as a graphic can be seen in Fig. 12. The total NIRS data collected is 720 data rows and 136 columns according to wavelength value from the sensor and visualized. as shown in Fig. 12. The wavelength used is between 1346.61 - 2556.24 nanometers (nm).



Fig. 9.   NeoSpectra development kit.



Fig. 10. NeoSpectra development kit with case and sample.

TABLE II.        EXAMPLE OF NIRS DATA

| Hour | Wavelength (nm) | | | | |
|------|---------|---------|-----|---------|---------|
| | *2556.24* | *2539.35* | *…* | *1351.35* | *1346.61* |
| 0 | 2.27 | 2.37 | … | 1.41 | 1.35 |
| 1 | 1.96 | 2.08 | … | 1.72 | 1.95 |
| 2 | 1.78 | 1.86 | … | 1.62 | 1.53 |
| 3 | 2.12 | 2.22 | … | 1.89 | 1.88 |
| 4 | 2.11 | 2.28 | … | 1.85 | 1.77 |
| 5 | 1.71 | 1.86 | … | 2.01 | 2.11 |
| 6 | 2.13 | 2.30 | … | 2.86 | 2.55 |
| 7 | 2.81 | 2.91 | … | 3.79 | 4.01 |



Fig. 12. Plot of original data spectrum.

## IV. PROPOSED FEATURE SELECTION MODEL

The proposed feature selection model aims to reduce the amount of data based on columns in the dataset, which by default has 136 columns. This feature selection concept is based on the line simplification model by reducing straight or slightly curved lines. The main stages of the feature selection process can be seen in Fig. 13.



Fig. 13. Flow process of proposed selection feature.

The feature selection process goes through four stages as follows :

### A. Mean and Single Data Line

In this process, 1 line is produced, which will represent all spectrum data. as shown in Fig. 14. Then, it is divided into one separate piece represented by one different color, as shown in Fig. 15. Then pair them with adjacent lines. as shown in Fig. 15. Then calculate the angle values of two adjacent lines as illustrated by Fig. 16.

### B. Iterative Line Simplification

The process of calculating all angles along a line produces 134 angle values. Then, the angle value data is sorted starting from the smallest value. Each iteration will eliminate one corner with the smallest value. The eliminated corners will correspond to the columns that will be eliminated as well. At this stage. a set of data columns is stored with the storage index in a sequence of iterations. The elimination stages are depicted in the diagram Fig. 17.



Fig. 14. Mean of all spectrums.



Fig. 15. Separated spectrum into single data line.



Fig. 16. Illustration of calculating the angle between two adjacent lines.



Fig. 17. Iterative angle and column elimination.

### C. Random Forest Regressor

Each iteration produces a new dataset with a reduction of 1 column of data. The new dataset will then enter the machine learning process to produce an $R^2$ value at each iteration. At this stage, the RFR parameters used are the default settings, as shown in Table III. Data splitting for the learning and testing process is 70% to 30%.

TABLE III.    DEFAULT PARAMETERS OF RFR [19]

| Parameter | Data type | Default value |
|---|---|---|
| n_estimators | | 100 |
| criterion | | squared_error |
| max_depth | | None |
| min_samples_split | | 2 |
| min_samples_leaf | | 1 |
| min_weight_fraction_leaf | float | 0.0 |
| max_features | int or float | 1.0 |
| max_leaf_nodes | int | None |
| min_impurity_decrease | float | 0.0 |
| bootstrap | bool | True |
| oob_score | bool or callable | False |
| n_jobs | int | None |
| random_state | int | None |
| verbose | int | 0 |
| warm_start | bool | False |
| ccp_alpha | non-negative float | 0.0 |
| max_samples | int or float | None |
| monotonic_cst | array-like of int of shape (n_features | None |

The results of $R^2$ at each iteration will be stored in an array. Then, it will find the highest $R^2$ value and location in the array to be used as the best column index set.

### D. The best result of the selected feature

The final stage is to determine the largest value from the collection of $R^2$ that has been accommodated in the array. To determine the highest $R^2$ value in this study, use the numpy library with the numpy.max() command [27]. To find out the iteration position of the highest $R^2$ value, also use the numpy library with the command numpy.argmax() [28]. This command will display the data index of the highest $R^2$ value.

To determine the best set of columns is to store the value of $R^2$ in each iteration. Then, find the highest value of the array to determine the index value. The index value is used to retrieve the set of columns.

## V. RESULTS

### A. Results from the Original Dataset and the Proposed Model

Of the six meat quality parameters, namely color, drip loss, pH, storage time, TCP, and water moisture, feature selection and machine learning models are applied alternately. Then, we will compare the prediction results using the original dataset with the dataset that has gone through the feature selection stage. The results of the comparison of $R^2$ values can be seen in IV.

Based on the $R^2$ value in Table IV, it can be seen that the line simplification-based feature selection model with the corner elimination method has succeeded in increasing the performance of the RFR algorithm in predicting all beef quality parameters. The increase in the $R^2$ value for all parameters can

be said to be satisfactory, with the smallest increase being in the predicted storage time parameter, which is only 0.028.

TABLE IV.    $R^2$ SCORE USING THE ORIGINAL DATASET AND USING FEATURE SELECTION

| Beef Quality Parameters | $R^2$ score | | |
|---|---|---|---|
| | *Original dataset* | *With the proposed feature selection* | *increment* |
| color | 0.597448 | 0.876992 | 0.279544 |
| drip loss | 0.891545 | 0.942723 | 0.051178 |
| pH | 0.796903 | 0.904225 | 0.107322 |
| storage time | 0.889070 | 0.916581 | 0.027511 |
| TPC | 0.720895 | 0.951388 | 0.230493 |
| water moisture | 0.539780 | 0.893200 | 0.353420 |

The results were very good, and the highest increase in the $R^2$ value was in the prediction of the water moisture parameter, namely 0.353. From all the increases in $R^2$ values, the average value can be taken to be 0.1749 or 17.49%.

### B. Results from the Original Dataset and SelectFromModel

As a comparison of performance results in this study, the feature selection model from Scikit-Learn, namely SelectFromModel. was also tested. The experimental results of implementing the SelectFromModel library can be seen in Table V.

TABLE V.    $R^2$ SCORE OF USING ORIGINAL DATASET AND SELECTFROMMODEL

| Beef Quality Parameters | $R^2$ score | | |
|---|---|---|---|
| | *Original dataset* | *SelectFromModel* | *increment* |
| color | 0.597448 | 0.835490 | 0.238042 |
| drip loss | 0.891545 | 0.917345 | 0.025800 |
| pH | 0.796903 | 0.846478 | 0.049575 |
| storage time | 0.889070 | 0.933890 | 0.044820 |
| TPC | 0.720895 | 0.928947 | 0.208052 |
| water moisture | 0.539780 | 0.823419 | 0.283639 |

Based on the $R^2$ value in Table V. it can be seen that feature selection by the SelectFromModel library can also improve the performance of the RFR algorithm. The smallest improvement was in the prediction of the drip loss parameter, namely 0.026, while the biggest improvement was in the prediction of the water moisture parameter, namely 0.2836. With the increase in the $R^2$ value in the prediction of all parameters, it can be said that feature selection using the SelectFromModel library works very well. Of all the increases in $R^2$. the average value can be taken to be 0.1417 or 14.17%.

## VI. DISCUSSION

### A. $R^2$ Score Comparison between the Proposed Model and SelectFromModel

In this stage. the results of the $R^2$ values from the proposed feature selection model and the SelectFromModel feature selection model are compared. The comparison results also contain the number of features selected based on their highest $R^2$ value, which can be seen in  Table VI.

TABLE VI. COMPARISON OF R2 SCORES BETWEEN THE PROPOSED MODEL AND SELECTFROMMODEL

| Beef Quality Parameters | Proposed Model | | SelectFromModel | |
|---|---|---|---|---|
| | $R^2$ | n feature | $R^2$ | n feature |
| color | 0.876992 | 84 | 0.835490 | 40 |
| drip loss | 0.942723 | 117 | 0.917345 | 20 |
| pH | 0.904225 | 57 | 0.846478 | 26 |
| storage time | 0.916581 | 71 | 0.933890 | 18 |
| TPC | 0.951388 | 56 | 0.928947 | 20 |
| water moisture | 0.893200 | 77 | 0.823419 | 31 |

Based on Table VI. it can be seen that the $R^2$ values of the proposed feature selection mode are higher than the SelectFromModel results, except for the $R^2$ value in the storage time parameter prediction. So overall, the average increase in the $R^2$ value of the proposed model is 17.49% higher than the average $R^2$ value of SelectFromModel. which produces 14.17%.

### B. Overview of Selected Features

The result of the feature selection model is a set of features in the form of data columns; in this study, the column names are wavelength values in nanometer (nm) units. The number of features produced by the proposed model and SelectFromModel is definitely less than the number of columns in the original dataset, so this feature selection also leads to a reduction in data dimensionality. For the differences in the number and features selected, a visualization was created for all parameters, which can be seen in Fig. 18 to Fig. 23.

Each color represents one feature selection model. The red color represents the mean of all NIRS data. The green color represents the mean of the data columns selected by the proposed model, while the blue represents the results from the SelectFromModel library.



Fig. 18. Spectrums overlay for the color parameter.



Fig. 19. Spectrum overlay for drip loss parameter.



Fig. 20. Spectrums overlay for pH parameter.

Fig. 21. Spectrums overlay for storage time parameter.



Fig. 22. Spectrums overlay for TPC parameter.



Fig. 23. Spectrums overlay for water moisture parameter.

## VII. CONCLUSION AND FUTURE WORK

This study has produced a feature selection model based on line simplification by eliminating angles in the average spectrum from beef NIRS data. The result of increasing RFR performance after using the proposed feature selection model is 17.49%. This result is higher than the average increase in R2 value produced by the SelectFromModel library of 14.17%. Apart from being able to increase prediction accuracy, this model can also reduce data dimensions, where fewer data will require a shorter time in the machine learning process.

For further work and development, the proposed feature selection model can be applied to deep learning algorithms. It can also be combined with RFR by applying hyperparameter tuning. The combination with hyperparameter tuning may require a longer time to find the solution set for the highest accuracy. However, RFR with hyperparameter tuning produces better accuracy compared to RFR with default parameters.

## REFERENCES

[1] M. Lima, R. Costa, I. Rodrigues, J. Lameiras, and G. Botelho, "A Narrative Review of Alternative Protein Sources: Highlights on Meat, Fish, Egg and Dairy Analogues," Foods, vol. 11, no. 14, p. 2053, Jul. 2022, doi: 10.3390/foods11142053.

[2] M. Molfetta et al., "Protein Sources Alternative to Meat: State of the Art and Involvement of Fermentation," Foods, vol. 11, no. 14, p. 2065, Jul. 2022, doi: 10.3390/foods11142065.

[3] B. Fletcher et al., "Advances in meat spoilage detection: A short focus on rapid methods and technologies," CyTA - J. Food, vol. 16, no. 1, pp. 1037–1044, Jan. 2018, doi: 10.1080/19476337.2018.1525432.

[4] V. Tesson, M. Federighi, E. Cummins, J. de Oliveira Mota, S. Guillou, and G. Boué, "A Systematic Review of Beef Meat Quantitative Microbial Risk Assessment Models," Int. J. Environ. Res. Public Health, vol. 17, no. 3, p. 688, Jan. 2020, doi: 10.3390/ijerph17030688.

[5] W. Barragán-Hernández, L. Mahecha-Ledesma, J. Angulo-Arizala, and M. Olivera-Angel, "Near-Infrared Spectroscopy as a Beef Quality Tool to Predict Consumer Acceptance," Foods, vol. 9, no. 8, p. 984, Jul. 2020, doi: 10.3390/foods9080984.

[6] A. Goi, J.-F. Hocquette, E. Pellattiero, and M. De Marchi, "Handheld near-infrared spectrometer allows on-line prediction of beef quality traits," Meat Sci., vol. 184, p. 108694, Feb. 2022, doi: 10.1016/j.meatsci.2021.108694.

[7]    A. Sahar et al., "Online Prediction of Physico-Chemical Quality Attributes of Beef Using Visible—Near-Infrared Spectroscopy and Chemometrics," Foods, vol. 8, no. 11, p. 525, Oct. 2019, doi: 10.3390/foods8110525.

[8]    N. Patel, H. Toledo-Alvarado, A. Cecchinato, and G. Bittante, "Predicting the Content of 20 Minerals in Beef by Different Portable Near-Infrared (NIR) Spectrometers," Foods, vol. 9, no. 10, p. 1389, Oct. 2020, doi: 10.3390/foods9101389.

[9]    S. Savoia et al., "Prediction of meat quality traits in the abattoir using portable and handheld near-infrared spectrometers," Meat Sci., vol. 161, p. 108017, Mar. 2020, doi: 10.1016/j.meatsci.2019.108017.

[10]   S. Savoia, A. Albera, A. Brugiapaglia, L. Di Stasio, A. Cecchinato, and G. Bittante, "Prediction of meat quality traits in the abattoir using portable near-infrared spectrometers: heritability of predicted traits and genetic correlations with laboratory-measured traits," J. Anim. Sci. Biotechnol., vol. 12, no. 1, p. 29, Dec. 2021, doi: 10.1186/s40104-021-00555-5.

[11]   I. M. N. Perez, L. J. P. Cruz-Tirado, A. T. Badaró, M. M. de Oliveira, and D. F. Barbin, "Present and future of portable/handheld near-infrared spectroscopy in chicken meat industry," NIR news, vol. 30, no. 5–6, pp. 26–29, Aug. 2019, doi: 10.1177/0960336019861476.

[12]   M. Simoni, A. Goi, M. De Marchi, and F. Righi, "The use of visible/near-infrared spectroscopy to predict fibre fractions, fibre-bound nitrogen and total-tract apparent nutrients digestibility in beef cattle diets and faeces," Ital. J. Anim. Sci., vol. 20, no. 1, pp. 814–825, Jan. 2021, doi: 10.1080/1828051X.2021.1924884.

[13]   C. N. Sánchez, M. T. Orvañanos-Guerrero, J. Domínguez-Soberanes, and Y. M. Álvarez-Cisneros, "Analysis of beef quality according to color changes using computer vision and white-box machine learning techniques," Heliyon, vol. 9, no. 7, p. e17976, Jul. 2023, doi: 10.1016/j.heliyon.2023.e17976.

[14]   T. Qiao, J. Ren, C. Craigie, J. Zabalza, C. Maltin, and S. Marshall, "Quantitative Prediction of Beef Quality Using Visible and NIR Spectroscopy with Large Data Samples Under Industry Conditions," J. Appl. Spectrosc., vol. 82, no. 1, pp. 137–144, Mar. 2015, doi: 10.1007/s10812-015-0076-1.

[15]   R. Kasarda, N. Moravčíková, G. Mészáros, M. Simčič, and D. Zaborski, "Classification of cattle breeds based on the random forest approach," Livest. Sci., vol. 267, p. 105143, Jan. 2023, doi: 10.1016/j.livsci.2022.105143.

[16]   Y. Lin, J. Ma, D.-W. Sun, J.-H. Cheng, and Q. Wang, "A pH-Responsive colourimetric sensor array based on machine learning for real-time monitoring of beef freshness," Food Control, vol. 150, p. 109729, Aug. 2023, doi: 10.1016/j.foodcont.2023.109729.

[17]   H. Liu, "Feature Selection," in Encyclopedia of Machine Learning, Boston, MA: Springer US, 2011, pp. 402–406. doi: 10.1007/978-0-387-30164-8_306.

[18]   J. Miao and L. Niu, "A Survey on Feature Selection," Procedia Comput. Sci., vol. 91, pp. 919–926, 2016, doi: 10.1016/j.procs.2016.07.111.

[19]   F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," J. Mach. Learn. Res., vol. 12, no. 85, pp. 2825–2830, 2011, [Online]. Available: http://jmlr.org/papers/v12/pedregosa11a.html.

[20]   L. Breiman, "Random forests," Mach. Learn., vol. 45, no. 1, 2001, doi: 10.1023/A:1010933404324.

[21]   J. Rogers and S. Gunn, "Identifying Feature Relevance Using a Random Forest BT - Subspace, Latent Structure and Feature Selection," 2006, pp. 173–184.

[22]   D. Chicco, M. J. Warrens, and G. Jurman, "The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE and RMSE in regression analysis evaluation," PeerJ Comput. Sci., vol. 7, p. e623, Jul. 2021, doi: 10.7717/peerj-cs.623.

[23]   K. B. Beć, J. Grabska, and C. W. Huck, "Near-Infrared Spectroscopy in Bio-Applications," Molecules, vol. 25, no. 12, p. 2948, Jun. 2020, doi: 10.3390/molecules25122948.

[24]   D. Tejerina, M. Oliván, S. García-Torres, D. Franco, and V. Sierra, "Use of Near-Infrared Spectroscopy to Discriminate DFD Beef and Predict Meat Quality Traits in Autochthonous Breeds," Foods, vol. 11, no. 20, p. 3274, Oct. 2022, doi: 10.3390/foods11203274.

[25]   M. Peyvasteh, A. Popov, A. Bykov, and I. Meglinski, "Meat freshness revealed by visible to near-infrared spectroscopy and principal component analysis," J. Phys. Commun., vol. 4, no. 9, p. 095011, Sep. 2020, doi: 10.1088/2399-6528/abb322.

[26]   "Si-Ware Systems announces NeoSpectra Micro Development Kits | NeoSpectra." https://www.si-ware.com/featured-resources/si-ware-systems-announces-neospectra-micro-development-kits (accessed Jun. 19, 2023).

[27]   "numpy.max — NumPy v1.26 Manual." https://numpy.org/doc/stable/reference/generated/numpy.max.html (accessed Nov. 16, 2023).

[28]   "numpy.argmax — NumPy v1.26 Manual." https://numpy.org/doc/stable/reference/generated/numpy.argmax.html (accessed Nov. 16, 2023).

# Research on Spatial Accessibility Measurement Algorithm for Sanya Tourist Attractions Based on Seasonal Factor Adjustment Analysis

Xiaodong Mao[1]*, Yan Zhuang[2]

School of Tourism and Health Industry, Sanya Institute of Technology, Sanya, Hainan, 572011, China[1]
University of Sanya, Sanya, Hainan, 572011, China[2]

*Abstract*—Seasonal factors will lead to changes in tourists' demand for scenic spots in different seasons, which will affect the traffic network and road conditions, and then affect the convenience and efficiency of tourists arriving at scenic spots. Based on the adjustment and analysis of seasonal factors, this study puts forward an algorithm for measuring the spatial accessibility of Sanya tourist attractions. Principal component analysis is used to denoise the data of Sanya tourist attractions in different seasons, and independent component analysis is used to extract the data characteristics of Sanya tourist attractions in different seasons after denoising. On this basis, the spatial accessibility index of Sanya tourist attractions is calculated by combining the spatial information of Sanya tourist attractions and GIS technology, and the spatial accessibility of Sanya tourist attractions is analyzed, and the spatial accessibility measurement model of Sanya tourist attractions is constructed to realize the spatial accessibility measurement of Sanya tourist attractions. The experimental results show that the spatial accessibility measurement method of Sanya tourist attractions is effective, which can effectively improve the accuracy of accessibility measurement and shorten the accessibility measurement time. It aims to help decision makers plan and optimize tourist routes and improve the efficiency and convenience of tourists arriving at their destinations.

*Keywords—Seasonal factors; adjustment analysis; Sanya Tourist Attractions; spatial accessibility measure; GIS technology*

## I. INTRODUCTION

As a famous seaside tourist city in China, Sanya has a unique natural environment and rich tourism resources. With the rapid development of tourism, the spatial accessibility of tourist attractions has become an important factor affecting the development of tourism. However, seasonal factors have a significant impact on the spatial accessibility of tourist attractions [1-3]. Therefore, it is of great practical significance to study the spatial accessibility measurement algorithm of Sanya tourist attractions based on seasonal factor adjustment analysis. Sanya has many famous tourist attractions, such as nansan, Yalong Bay Tropical Paradise Forest Park and Tianya Haijiao. With the development of tourism, the spatial accessibility of tourist attractions has become an important indicator to measure the development level of tourism in a region [4]. Tourist attractions with high spatial accessibility can attract more tourists, create more economic benefits and promote the sustainable development of tourism in Sanya.

Zhou Haitao [5] and others revealed the spatial distribution characteristics of red tourist attractions in Inner Mongolia by means of kernel density and geographical concentration index. Based on the road planning function of Gaode map, the road traffic conditions were obtained in real time, and the spatial accessibility measurement model of red tourist attractions was constructed. The influencing factors of accessibility differences were clarified by using geographical detectors. This method can comprehensively consider the spatial distribution characteristics of scenic spots, traffic network and road traffic conditions, and comprehensively evaluate the accessibility of scenic spots from multiple angles. However, depending on the path planning function and road traffic data provided by Gaode map may have a certain impact on the measurement results. Wang Hao [6] and others, based on the statistical yearbook of Xinjiang road network and the official data of Xinjiang Culture and Tourism Department, analyzed the spatial distribution characteristics of scenic spots of 4A level and above in Xinjiang by using the network analysis method, kernel density analysis method and geographical concentration index of GIS, and analyzed the spatial accessibility from two aspects: scenic accessibility and regional accessibility. This method can accurately calculate the shortest path and time between scenic spots, thus providing more accurate results of scenic accessibility analysis. However, the reliability of the data needs to be verified. Wei Liu [7] and others use the spatial grid method to divide the research scope, count the number of residents in the research scope, and use the API interface of the network map platform to obtain the expected travel time data of residents. On this basis, according to the accessibility model, the fairness of public transport accessibility in Xi 'an is analyzed by using fairness coefficient and Lorenz curve. This method can quickly obtain the calculation results, but the data source is limited. Ma Shuhong [8] and others calculate the travel cost of public transport according to the data planned by real-time routes. Arcgis is used to analyze the accessibility of urban agglomerations, and theil index and fairness coefficient are used to get the difference characteristics of public transport accessibility. The calculation results of accessibility and fairness of this method are accurate and the coverage is limited.

Based on this, this paper puts forward an algorithm for measuring the spatial accessibility of Sanya tourist attractions based on the adjustment and analysis of seasonal factors, aiming at exploring how to better eliminate the influence of seasonal factors, improve tourists' travel experience and

promote the sustainable development of Sanya tourism through in-depth research on the spatial accessibility of Sanya tourist attractions. By extracting the data characteristics of Sanya tourist attractions in different seasons, this paper analyzes the spatial accessibility of Sanya tourist attractions, constructs a spatial accessibility measurement model of Sanya tourist attractions, and effectively adjusts and analyzes the spatial accessibility of Sanya tourist attractions due to seasonal factors. This method can effectively measure the spatial accessibility of Sanya tourist attractions, improve the measurement accuracy and shorten the measurement time. By studying the influencing factors of the spatial accessibility of Sanya tourist attractions, it provides scientific basis and decision support for the prosperity and development of Sanya tourism. The research contribution of this paper:

*1)* Denoising the data of Sanya tourist attractions in different seasons, extracting the data characteristics of Sanya tourist attractions in different seasons, considering the changes in seasonal demand, traffic network, and other factors, assisting decision makers in planning and optimizing tourist routes, improving the efficiency and convenience of tourists arriving at their destinations.

*2)* Based on the spatial information of Sanya tourist attractions and GIS technology, the spatial accessibility index of Sanya tourist attractions is calculated, so as to evaluate the spatial accessibility of Sanya tourist attractions, provide scientific basis and decision support for the development and planning of Sanya tourism, and make positive contributions to improving tourism efficiency and promoting economic development.

## II. Feature Extraction of Sanya Tourist Attractions Data in Different Seasons

Due to the fact that different seasons attract different types of tourists and have an impact on the number of tourists, transportation network, etc. By extracting data features of Sanya tourist attractions in different seasons, it is possible to analyze the changes in demand for Sanya tourist attractions and the operation of transportation and facilities during different seasons. Count the number of tourists to tourist attractions in Sanya during different seasons, understand the characteristics and travel preferences of the tourist group, help formulate more accurate Sanya tourism planning and marketing strategies, conduct seasonal factor adjustment analysis, optimize traffic management and resource allocation, and improve the spatial accessibility of Sanya tourist attractions. Therefore, in order to effectively measure the spatial accessibility of Sanya tourist attractions, first, principal component analysis is used to denoise the data of Sanya tourist attractions in different seasons. Then, independent component analysis is used to extract the features of the denoised data of Sanya tourist attractions in different seasons.

### A. Sanya Tourist Attractions Data Noise Reduction in Different Seasons

Based on the principle of constrained optimization, obtain abnormal data of tourist attractions in Sanya in different seasons and flexibly represent them.

The data of tourist attractions in Sanya for different seasons has $n$ dimensions, and each dimension includes a layer $d_{rj}$, which also has $l_r$ dimension layers. Set the value at position $i$ of the $r$ dimension $j$ layer in the data of Sanya tourist attractions in different seasons to $y_{i_1,i_2,\ldots,i_{rj}}$. The process of obtaining the expected value $\hat{y}_{i_1,i_2,\ldots,i_{rj}}$ of Sanya tourist attraction data in different seasons is as follows:

$$\hat{y}_{i_1,i_2,\ldots,i_{rj}} = y_{i_1,i_2,\ldots,i_{rj}}\left(\gamma_{i_{rj}|d_{rj}}^{G} \mid G \subset \left\{d_{rj} \mid 1 \le r \le n, 1 \le j \le l_r\right\}\right) \quad (1)$$

In Formula (1), $\gamma_{i_{rj}|d_{rj}}$ is the expected coefficient of Sanya tourist attraction data in different seasons, and $G$ is the constraint condition in Sanya tourist attraction data in different seasons. The abnormal data $s_{i_1,i_2,\cdots,i_{rj}}$ of tourist attractions in Sanya obtained from different seasons is as follows:

$$s_{i_1,i_2,\cdots,i_{rj}} = \frac{\left| y_{i_1,i_2,\ldots,i_{rj}} - \hat{y}_{i_1,i_2,\ldots,i_{rj}} \right|}{\sigma_{i_1,i_2,\cdots,i_{rj}}} \quad (2)$$

In Formula (2), $\sigma_{i_1,i_2,\cdots,i_{rj}}$ is the abnormal data index of Sanya tourist attractions in different seasons.

Based on the above calculation results, the abnormal data of Sanya tourist attractions in different seasons will be removed, and then the principal component analysis method [9-11] will be used to denoise the data in Sanya tourist attractions in different seasons. Principal component analysis (PCA) is a commonly used data denoising method, which transforms the original data into a new coordinate system by linear transformation, so that the largest variance appears on the first coordinate axis, the second largest variance appears on the second coordinate axis, and so on. The final principal components are linear combinations of original data changes, and these principal components are sorted according to their variance.

Construct a three-dimensional coordinate system using the component decomposition method in Sanya tourist attractions of different seasons, and the data vectors in Sanya tourist attractions of different seasons in the coordinates are:

$$Z(x) = \left[Z_1(x), Z_2(x), \cdots Z_p(x)\right]^T = S_p(x) + N_p(x) \quad (3)$$

In Formula (3), $S_p(x), N_p(x)$ represents the $p$-dimensional signal and noise vector of Sanya tourist attraction data in different seasons.

By using the minimum noise separation transformation method, the data from Sanya tourist attractions in different seasons are linearly transformed, and the process is as follows:

$$Y_i(x) = a_i^T Z(x), i = 1, 2, \cdots, p \tag{4}$$

In Formula (4), $a_i^T$ is the transformation coefficient of data from Sanya tourist attractions in different seasons, and $Y_i(x)$ is the transformed data from Sanya tourist attractions in different seasons. Through the above transformation, it can be seen that each band in $Y_i(x)$ exists independently of each other, and the obtained result is the maximum signal-to-noise ratio $SNR_{Y_i}$ of data from Sanya tourist attractions in different seasons. The acquisition process is as follows:

$$SNR_{Y_i} = \frac{\mathrm{var}\, Y_i(x)\{a_i^T S_p(x)\} V}{\mathrm{var}\, Y_i(x)\{a_i^T N_p(x)\} B} \tag{5}$$

In Formula (5), $V, B$ is the covariance matrix of signal and noise in the data of Sanya tourist attractions in different seasons. Based on the above results, complete the data denoising processing for Sanya tourist attractions in different seasons.

*B. Feature Extraction of Sanya Tourist Attractions Data in Different Seasons*

Independent component analysis (ICA) is a data analysis method applied in signal processing, neural network and other fields. Its core idea is to assume that the observed multidimensional signal is a linear mixture of multiple independent components, each of which has its own statistical distribution. The purpose of ICA is to decompose the mixed signal into independent components through certain linear transformation, so as to realize signal separation and feature extraction. When ICA is used to extract the data features of Sanya tourist attractions in different seasons, it is generally necessary to preprocess the data first, such as using PCA to denoise the data to eliminate the noise and interference information in the data. Then, ICA algorithm is used to linearly transform the preprocessed data to extract the independent components hidden in the data. Using independent component analysis method, extract the features of Sanya tourist attractions data in different seasons after denoising.

Set $Z(x)$ as the original data vector of Sanya tourist attractions in different seasons, and the data vector of Sanya tourist attractions in different seasons after minimum noise separation transformation is $Y(x)$. Retain the first $k$ term in $Y(x)$ and reset the other $p - k$ terms to zero, in order to obtain the dataset of Sanya tourist attractions in different

seasons for the first $k$ term of $Y(x)$ and perform minimum noise separation inverse transformation on it. The process is as follows:

$$\begin{cases} Y^{(k)}(x) = \left[ Y_1(x), Y_2(x), \cdots, Y_k(x), 0, \cdots, 0 \right]^T \\ Q(x) = \left[ 1Q(x), 2Q(x), \cdots, pQ(x) \right]^T \\ Q_i(x) = SNR_{Y_i} \sum_{t=1}^{p} a_t Y_t \end{cases} \tag{6}$$

In Formula (6), $Y^{(k)}(x)$ is the dataset of Sanya tourist attractions in different seasons obtained from the first $k$ term of $Y(x)$, $Q(x)$ is the dataset of Sanya tourist attractions in different seasons obtained from the minimum noise separation inverse transformation, $a_t$ is the conversion amount during the transformation, and $p$ is a constant. Based on the above calculation results, construct an orthogonal matrix for the dataset of Sanya tourist attractions in different seasons, and the process is as follows:

$$W = Q_i(x) A = \left[ 1Q(x), 2Q(x), \cdots, pQ(x) \right]^T \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1p} \\ a_{21} & a_{22} & \cdots & a_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \cdots & a_{pp} \end{bmatrix} \tag{7}$$

In Formula (7), $A$ is a singular matrix with zero eigenvalues. Calculate the matrix based on the third-order center distance to obtain the independent component skewness of Sanya tourist attraction data for different seasons. The process is as follows:

$$\det(W) = \det(0)° \det(A) = 0 \tag{8}$$

In Formula (8), $\det(A)$ is the independent component kurtosis of the obtained data of Sanya tourist attractions in different seasons. Sort the data of Sanya tourist attractions in different seasons based on the obtained independent component skewness, in order to obtain the characteristics of Sanya tourist attraction data in different seasons.

III. MEASUREMENT OF SPATIAL ACCESSIBILITY OF TOURIST ATTRACTIONS IN SANYA

The measurement of spatial accessibility of Sanya tourist attractions is based on the data characteristics of Sanya tourist attractions in different seasons, combined with the spatial information of Sanya tourist attractions and GIS technology, to calculate the spatial accessibility index of Sanya tourist attractions, and evaluate the spatial accessibility of Sanya tourist attractions. These results can provide scientific basis for tourism planners in Sanya, helping them understand the changes in spatial accessibility of tourist attractions in Sanya, and carry out reasonable resource development and Sanya tourism route design.

*A. Analysis of Spatial Accessibility of Tourist Attractions in Sanya*

Assuming $SY_m$ represents the temporal accessibility of the spatial node $m$ of Sanya tourist attractions, it describes the shortest time average value of the spatial node $m$ of Sanya tourist attractions reaching other nodes. The mean temporal accessibility of all nodes is the temporal accessibility $SY$ of the entire Sanya tourist attraction space, and its calculation formula is as follows:

$$\begin{cases} SY_m^s = \det(W) \sum_{m=1,v\neq 1}^{\omega} SY_{mv}^s / (Y_{mv}^s - 1), m \in [1, \omega] \\ SY^s = \sum_{m=1}^{\omega} SY_m^s / Y_{mv}^s \end{cases}$$
(9)

In Formula (9), $\omega$ represents the number of nodes in the Sanya region, and $SY_m^s$ represents the time accessibility of node $m$ under travel mode $s$, $SY_{mv}^s$ represents the time barrier between node $m$ and node $v$ in the Sanya tourist attraction space using travel mode $s$, and $SY^s$ represents the time accessibility corresponding to travel mode $s$ in the study area.

Construct a set of travel modes $S$, which includes three types of travel modes: car travel, motorcycle travel, and bicycle travel. Assuming $Y_{mv}^s$ is the travel time generated by travel mode $s$ on the designated road section, it can be calculated using Formula (10):

$$Y_{mv}^s = \frac{r_{mv}}{b_{mv}^s}$$
(10)

In Formula (10), $r_{mv}$ is the distance between node $m$ and node $v$ in the space of Sanya tourist attractions, and $b_{mv}^s$ is the speed of the existing section $r_{mv}$ under travel mode $s$.

Assuming that $SV_m$ represents the economic cost distance corresponding to node $m$ in the Sanya tourist attraction space, it describes the minimum average cost of node $m$ to other nodes in the Sanya tourist attraction space. The average time accessibility corresponding to all nodes is the time accessibility $SV$ corresponding to the Sanya tourist attraction space, and its calculation formula is as follows:

$$\begin{cases} SV_m^s = \det(W) \sum_{m=1,v\neq 1}^{\omega} V_{mv}^s / (Y_{mv}^s - 1) \\ SV^s = \sum_{m=1}^{\omega} SV_m^s / Y_{mv}^s \end{cases}$$
(11)

In Formula (11), $SV_m^s$ represents the cost accessibility of node $m$ under travel mode $s$, $V_{mv}^s$ represents the cost of node $m$ using travel mode $s$ to reach node $v$ in the Sanya tourist attraction space, and $SV^s$ represents the cost accessibility of using travel mode $s$ in the study area.

The generalized travel cost usually includes the following types:

*1)* Parameter comprehensive cost $V_{mv}^z$:

$$V_{mv}^z = Y_{mv}^s + VOT \times \beta Y_{mv}^s \left(\frac{Y_{mv}^s}{Y_c}\right)^{\chi}$$
(12)

In Formula (12), $VOT$ is the unit time cost, $\beta$ and $\chi$ are adjustment coefficients, and $Y_c$ is the time threshold.

*2)* Parameter time cost $V_{mv}^S$:

$$V_{mv}^S = VOTY_{mv}^s + VOTY_{mv}^s \beta \left(\frac{Y_{mv}^s}{Y_c}\right)^{\chi}$$
(13)

*3)* Comprehensive cost $V_{mv}^{ZH}$:

$$V_{mv}^{ZH} = VOTY_{mv}^s + Sin_{mv}^s$$
(14)

In Formula (14), $Sin_{mv}^s$ is the travel cost incurred from node $m$ using travel method $s$ to node $v$.

*4)* Travel cost $V_{mv}^{CX}$:

$$V_{mv}^{CX} = \frac{VOTr_{mv}}{b_{mv}^s} + r_{mv}\omega_{mv}^s$$
(15)

In Formula (15), $\omega_{mv}^s$ is the fuel consumption generated by using travel mode $s$ on section $r_{mv}$.

Using the weighted shortest travel time method, calculate the time of the Sanya tourist attraction space at each node to reflect the spatial accessibility of Sanya tourist attractions. The calculation formula is:

$$\begin{cases} A_x = \sum_{y=1}^{\tau} \lambda_y \delta_{xy} / \sum_{y=1}^{\tau} \lambda_y \\ \lambda_y = \sqrt{p_x e_x} \end{cases}$$
(16)

In Formula (16), $A_x$ represents the weighted time for the shortest travel in Sanya's tourist attraction space, and $\delta_{xy}$ represents the minimum time that should be consumed for the shortest distance in Sanya's tourist attraction space, $\lambda_y$

represents the capacity of Sanya's tourist attraction space and the degree of connection with other cities, $p_x$ represents the number of people in Sanya's tourist attraction space, $e_x$ represents the gross economic product of Sanya's tourist attraction space, $y$ represents the location parameter, and $\tau$ represents the spatial range of Sanya's tourist attraction. The smaller the $A_x$ value, the higher the spatial accessibility of tourist attractions in Sanya is.

The comprehensive weighted average travel time is used to measure the accessibility of tourist attractions in Sanya, and the calculation formula is:

$$IA_x = \sum A_{xx_i} r_{x_i} \tag{17}$$

In Formula (17), $IA_x$ represents the spatially weighted shortest average travel time of tourist attractions in Sanya, $A_{xx_i}$ represents the shortest average travel time, and $r_{x_i}$ represents the travel weight.

Based on the above calculation formula and the urban spatial gravity model, the shortest spatial time of tourist attractions in Sanya is calculated using the following formula:

$$I_{mn} = IA_x \mu_m \mu_n / d_{mn}^b \tag{18}$$

In Formula (18), $I_{mn}$ represents the spatial interaction force between tourist attractions in Sanya, $d_{mn}^b$ represents the shortest time corresponding to the shortest spatial distance of tourist attractions in Sanya, $\mu_m$ and $\mu_n$ represent the spatial scale of tourist attractions in Sanya, and $b$ represents the spatial distance friction index of tourist attractions in Sanya.

Using the potential model, calculate the spatial interaction potential $J_x$ of tourist attractions in Sanya, and the calculation formula is:

$$J_x = \sum_{m=1}^{n} I_{mn} + \mu_m \mu_n / d_{mn}^b \tag{19}$$

Traffic impedance is related to factors such as residents' travel modes and transportation environment, and is a physical factor that measures the difficulty of travel. Combining traffic impedance with the spatial road characteristics of Sanya tourist attractions, taking into account both time and cost impedance factors, this paper introduces them into the spatial accessibility analysis of Sanya tourist attractions, and improves the spatial accessibility analysis of Sanya tourist attractions. The function expression is:

$$\xi_i = \frac{(\theta_1 T_i + \theta_2 F_i)}{w_j} J_x \tag{20}$$

In Formula (20), $\xi_i$ represents the traffic resistance value, $\theta_1$ and $\theta_2$ represent the coefficients to be labeled, $T_i$ represents the time resistance value, $F_i$ represents the cost resistance value, and $w_j$ represents the regional importance.

*B. Construction of Spatial Accessibility Measurement Model for Sanya Tourist Attractions*

Based on the analysis of spatial accessibility of tourist attractions in Sanya, a measurement model for spatial accessibility of tourist attractions in Sanya is constructed. According to the time schedule and various data obtained, there are four main types of models for measuring the spatial accessibility of tourist attractions in Sanya, namely physical accessibility evaluation models, comprehensive index models based on residents' spatial preferences for tourist attractions in Sanya, analysis models based on differences in residents' spatial motivations for choosing tourist attractions in Sanya, and analysis models based on temporal and spatial behavior, use these four models to measure the spatial accessibility of tourist attractions in Sanya.

*1) Physical accessibility evaluation model:* This type of model is often combined with shortest distance, coverage, gravity method, etc. to describe the spatial distribution of transportation vehicles in Sanya's tourist attractions. The most commonly used indicators for evaluating the supply and demand of transportation vehicles are the distance from residents' residential areas to the location of tourist attractions in Sanya, as well as the spatial scale, capacity, and construction quality of tourist attractions in Sanya.

Based on the two evaluation indicators of shortest distance and scale capacity, the calculation formula is:

$$\begin{cases} D_m = \min_{v} d_{mv} \\ C_m = \dfrac{\xi_i \sum S_v}{P_m} \end{cases} \tag{21}$$

In Formula (21), $D_m$ represents the shortest distance from the starting point to the destination, and $C_m$ represents the spatial capacity of tourist attractions in Sanya, $d_{mv}$ represents the geometric centroid distance, $\min$ represents the minimum function value, $S_v$ represents the total area of spatial facilities construction in Sanya tourist attractions, and $P_m$ represents the total population of the city.

Due to the fact that physical evaluation models are more commonly applied to the same transportation vehicle and have certain limitations, the two evaluation indicators of shortest distance and scale capacity have been modified. The shortest distance from residents' residences to the location of Sanya tourist attractions is $D_{mm}$, and the effective coverage range of

Sanya tourist attractions is $C_{vv}$. The physical accessibility evaluation model is:

$$R = D_{mm} C_{vv} \xi_i \sum_{d_{mv} < K} E\rho$$

(22)

In Formula (22), $K$ represents the effective coverage area, $E$ represents the accessibility difference index for residents to access the spatial facilities of Sanya tourist attractions, and $\rho$ represents the population density of transportation users.

According to geographic information system (GIS) classification [12-13], differences in residents' age, income, and other factors are statistically analyzed. Combined with the distribution of population living space, the degree to which different residential areas are covered by the spatial facilities of Sanya tourist attractions is compared, and an evaluation of the physical accessibility of Sanya tourist attractions is completed.

*2) Comprehensive index model Based on residents' spatial preference for Sanya tourist attractions:* For residents, each type of Sanya tourist attraction spatial facility has its unique characteristics, which can meet the specific needs of specific subjects. At the same time, residents have various preferences for Sanya tourist attraction spatial facilities. Due to the fact that the audience level of tourist attractions in Sanya depends on multiple reasons, a comprehensive index model (IEI) is established using two indicators: attractiveness and separation to analyze residents' demand and preference for Sanya tourist attractions. The comprehensive index model takes residents' preference for Sanya tourist attractions as a weight coefficient and uses Likert scale scale to represent the spatial preference of Sanya tourist attractions. The expression of the comprehensive index model for Sanya tourist attractions is:

$$F_{ij(k)} = P_k \times W_{j(k)} \times s_{ij}^{-\alpha}$$

(23)

In Formula (23), $P_k$ represents the preference coefficient of different residents for the spatial facilities of different Sanya tourist attractions, $W_{j(k)}$ represents the indicator that attracts residents to choose, $s$ represents the measurement value of the spatial separation degree of Sanya tourist attractions, and $\alpha$ represents the spatial separation coefficient of Sanya tourist attractions.

Based on GIS output classification, a visual map with $F_{ij(k)}$ values is obtained to analyze the spatial distribution of tourist attractions in Sanya.

*3) Analysis model Based on spatial motivation differences of residents choosing Sanya tourist attractions:* This type of analysis model is based on the gravity model and presents the analysis results through GIS and spatial distribution models. Considering the behavioral ability of residents to choose spatial facilities in Sanya tourist attractions and the differences in simulated space, the selection motivation is divided into

$\beta_k$ and $\alpha_k$. The calculation formula for the analysis model of spatial motivation differences in Sanya tourist attractions is:

$$G_{ij(r)} = \frac{M_{j(r)} \times D_{j(r)}}{\alpha_k \times \beta_k \times T_{kj}}$$

(24)

In Formula (24), $T_{kj}$ represents the spatial value of Sanya tourist attractions, $r$ represents the serial number of Sanya tourist attractions spatial facilities, $M_{j(r)}$ represents the scale of Sanya tourist attractions spatial facilities, and $D_{j(r)}$ represents the service range of Sanya tourist attractions spatial facilities.

*4) Analysis model based on temporal and spatial behavior:* This model mainly analyzes the actual travel schedules of residents and the spatial travel schedules of Sanya tourist attractions, and combines them with spatiotemporal variables to calculate the comprehensive service deprivation coefficient of Sanya tourist attractions. The model considers the matching degree between residents' personal schedules and the spatial schedules of Sanya tourist attractions. Analyze the accessibility of this model from three perspectives:

*a)* The calculation formula for measuring the spatiotemporal accessibility of tourist attractions in Sanya is:

$$\begin{cases} U_{AB}(z, t_1, t_2) = \max(z, t_1, t_2, f) \\ U_{T_\varepsilon T_\varepsilon}(z, t_1, t_2) = \min(z, f) U_{AB}(z, t_1, t_2) \end{cases}$$

(25)

In Formula (25), $t_1$ and $t_2$ represent the comprehensive time period, $z$ represents the location of spatial facilities for tourist attractions in Sanya, and $f$ represents the number of spatial facilities for tourist attractions in Sanya during opening hours, $U_{AB}(z, t_1, t_2)$ represents the accessibility ability of residents to reach $f$ in both time and airspace during the comprehensive time period, while $U_{T_\varepsilon T_\varepsilon}(z, t_1, t_2)$ represents the accessibility ability of residents to reach $f$ in the shortest travel time during the comprehensive time period.

*b)* The calculation formula for measuring the spatiotemporal needs of spatial facilities in tourist attractions in Sanya is:

$$H_t(z, t_1, t_2) = \int_{t_1}^{t_2} N_t(z) dz / a$$

(26)

In Formula (26), $N_t(z)$ represents the individual residents' demand for the space of tourist attractions in Sanya, and $a$ represents the demand cycle.

Combining the space accessibility of Sanya tourist attractions with the space requirements of Sanya tourist

attractions, the analysis model based on time and space behavior is:

$$\eta(z,t_1,t_2) = \varpi \times \left[ U_{T_\varepsilon T_\varepsilon}(z,t_1,t_2) \times H_t(z,t_1,t_2) \right] \quad (27)$$

In Formula (27), $\varpi$ represents the weights of time and spatial indicators.

By calculating the spatiotemporal accessibility of Sanya tourist attractions and the spatiotemporal needs of facilities in the model, combined with temporal and spatial indicators, the results are output through GIS.

In order to facilitate calculation and analysis by unifying various indicators, the range standardization method is used to obtain the calculation formula as follows:

$$\varpi_j = (\varpi' - \min \varpi') / (\max \varpi' - \min \varpi') \quad (28)$$

In Formula (28), $\varpi_j$ represents the standardized values of time-domain and spatial indicators, $\varpi'$ represents the original values of time-domain and spatial indicators, $\max \varpi'$ represent the maximum function, and $\min \varpi'$ represents the minimum function.

Due to the strong objectivity and high numerical accuracy of the entropy method, it can effectively reduce the weight of indicators and be affected by subjective factors. Therefore, using the entropy method [14-15] to determine the weight of each indicator can better reflect the impact of evaluation indicators on the evaluation results. The calculation formula is:

$$\begin{cases} R_j = \varpi' / \sum_{j=1}^{n} \varpi_j \\ e_j = \sum_{j=1}^{n} (R_j \times \ln R_j) \\ g_j = 1 - e_j \\ W_j = \sum_{j=1}^{n} g_j \end{cases} \quad (29)$$

In Formula (29), $R_j$ represents the proportion of indicators, $e_j$ represents entropy, $g_j$ represents the coefficient of difference, and $W_j$ represents the weight coefficient.

The comprehensive accessibility measurement results of the four models are:

$$Z = \sum_{j=1}^{n} W_j \left[ G_{ij(r)} F_{ij(k)} C_{vv} D_{mm} \eta(z,t_1,t_2) \right] \quad (30)$$

Through the above steps, the spatial accessibility measurement of tourist attractions in Sanya can be achieved.

## IV. EXPERIMENTAL ANALYSIS

### A. Experimental Environment and Data Sources

As a tourist destination in the tropics, Sanya has many natural landscapes, such as beautiful beaches, magnificent mountains and tropical rainforests. These natural landscapes not only attract many tourists, but also become an important feature of Sanya tourism. In order to verify the validity of the spatial accessibility measurement algorithm of Sanya tourist attractions based on seasonal factor adjustment analysis, the experiment is implemented in MATLAB simulation software environment. Taking Coconut Dream Corridor, Yalong Bay, Haitang Bay, Jiajing Island, wuzhizhou Tourist Scenic Area, Fenjiezhou Island and Nanwan Monkey Island as spatial nodes, the data set of Sanya tourist attractions is divided. Through questionnaires, interviews and other means to collect tourists' evaluation of scenic spots, tourism motivation, consumption behavior and other information. Collect relevant data from government departments, tourism websites, scenic spots management offices and other channels, such as the number of tourists, tourism income, ticket sales, etc. Using big data technology, through analyzing data sources such as online search, social media and travel booking websites, we can get information about tourists' attention and preferences to Sanya tourist attractions.

There may be some noises and missing values in the collected tourism data, which need to be preprocessed, including data cleaning, missing value processing, abnormal value processing and so on. For example, removing duplicate data, correcting erroneous data, filling in missing values, etc. Due to the seasonal factors in Sanya, the popularity and demand of scenic spots may be different in different seasons. Summer is the peak of tourism in Sanya, and many beaches and water sports attractions will become more popular. In winter, some hot springs and tropical botanical gardens may be favored by tourists. Therefore, it is necessary to consider the change of seasonal demand when analyzing the data of Sanya tourist attractions. According to the seasonal characteristics of Sanya tourist attractions, the data can be divided into different seasonal segments. Generally speaking, the tourist season in Sanya is mainly concentrated in winter and summer vacations and holidays, and the data can be divided into seasons according to the date. Based on the data collection, preprocessing and seasonal segmentation of Sanya tourist attractions, the accessibility of Sanya tourist attractions at each node is measured by using the spatial calculation method of Sanya tourist attractions.

### B. Analysis of the Measurement Effect of Spatial Accessibility of Tourist Attractions in Sanya

Accessibility can reflect the difficulty of interconnecting one region with other regions. In order to verify the effectiveness of the spatial accessibility measurement method for tourist attractions in Sanya, the accessibility of this study is understood as the average travel time from a tourist attraction in a certain administrative region of Sanya City to tourist attractions in other administrative regions. The average travel time $s_\tau$ of tourist attraction $\tau$ in Sanya is expressed as:

$$s_\tau = \sum_{\xi=1,\zeta=1}^{\psi} g_{\xi\zeta} / \psi$$

$$(31)$$

In Formula (31), $\psi$ is the number of tourist attractions in Sanya, and $g_{\xi\zeta}$ is the travel time from point $\xi$ in the region to point $\zeta$ through the shortest route in the transportation network.

Select Coconut Dream Corridor, Yalong Bay, Haitang Bay, Jiajing Island, wuzhizhou Tourist Scenic Area, Fenjiezhou Island and Nanwan Monkey Island as the evaluation index. The smaller the value, the better the spatial accessibility measurement effect of Sanya tourist attractions. The average travel time of tourist attractions in Sanya is shown in Table I.

TABLE I.        AVERAGE TRAVEL TIME OF TOURIST ATTRACTIONS IN SANYA

| Sanya Tourist Attractions | $S_\tau$ /min |
|---|---|
| Coconut Dream Corridor | 151.35 |
| Yalong Bay | 120.86 |
| Haitang Bay | 101.65 |
| Jiajing Island | 125.42 |
| Wuzhizhou Tourist Scenic Area | 78.24 |
| Fenjiezhou Island | 81.24 |
| Nanwan Monkey Island Ecotourism Area | 74.43 |

According to the data in Table I, the average travel time of seven tourist attractions in Sanya is 74.43min-151.35min. Among them, the Nanwan Monkey Island Ecotourism Area, with the shortest average travel time of seven tourist attractions in Sanya, is about 74.43min, and this tourist attraction in Sanya has the best spatial accessibility. Followed by Wuzhizhou tourist Scenic Area, about 78.24min, the average travel time of Sanya tourist attractions is the Coconut Dream Corridor, about 151.35min, and the spatial accessibility of Sanya tourist attractions is the worst. Therefore, the method in this paper can effectively measure the spatial accessibility of Sanya tourist attractions, and has a good effect on measuring the spatial accessibility of Sanya tourist attractions.

### C. Precision Analysis of Spatial Accessibility Measurement for Sanya Tourist Attractions

In the process of calculating the spatial accessibility of tourist attractions in Sanya, the Lorentz curve and fairness coefficient are key indicators for judging fairness. Among them, the fairness coefficient is based on the Lorentz curve, and the fairness coefficient value can be calculated based on the area ratio in the Lorentz curve graph. Therefore, the fairness coefficient is used to calculate the fairness level of spatial accessibility of tourist attractions in Sanya. Randomly select 10 time periods as experimental data samples, calculate the fairness coefficients for 10 time periods, and compare the fairness calculation results of the method of reference [7], the method of reference [8], and this article with the actual results. The comparison results of fairness coefficients between different methods are shown in Fig. 1.



Fig. 1.    Comparison results of fairness coefficients of different methods.

According to Fig. 1, the fairness coefficients calculated by different methods in different time periods are different. Comparing the actual results with the methods in this paper, literature [7] and literature [8], we can see that the methods in this paper are in good agreement with the actual results, while the methods in literature [7] and literature [8] are far from the actual results. This shows that the fairness calculation of this method is more accurate and can effectively improve the accuracy of spatial accessibility measurement of Sanya tourist attractions. This is because this method combines the advantages of principal component analysis and independent component analysis, and can consider the characteristics and problems of data more comprehensively. Principal component analysis can reduce the dimension of multivariate to a few principal components and extract the main features of data. Independent component analysis can extract the independent components from the data, further denoising and simplifying the data. By using these two methods comprehensively, the characteristics and problems of data can be considered more comprehensively, thus improving the accuracy of spatial accessibility measurement of tourist attractions.

### D. Time Analysis of Spatial Accessibility Measurement for Tourist Attractions in Sanya

On this basis, further validate the spatial accessibility measurement time of Sanya tourist attractions using the method proposed in this paper. Compare the method of reference [7] and the method of reference [8] with the method proposed in this paper, and obtain the spatial accessibility measurement time of Sanya tourist attractions using different methods, as shown in Fig. 2.

According to Fig. 2, with the increase of the number of time periods, the time for measuring the spatial accessibility of Sanya tourist attractions by different methods increases. When the time period reaches 10, the measurement time of the method in reference [7] is 14.7ms, that of the method in reference [8] is 17.6ms, and that of this method is 4.8 ms. Therefore, it takes a short time to measure the spatial accessibility of Sanya tourist attractions by this method. This is because, this method constructs physical accessibility

evaluation model, comprehensive index model based on residents' spatial preference for Sanya tourist attractions, analysis model based on residents' spatial motivation differences in choosing Sanya tourist attractions, and analysis model based on time domain and spatial behavior. The physical accessibility evaluation model can transform complex spatial relations into simple distance and time calculations; The comprehensive index model based on residents' spatial preference for Sanya tourist attractions and the analysis model based on residents' spatial motivation differences in choosing Sanya tourist attractions can measure tourists' preferences and choices from different angles, thus avoiding complex spatial optimization algorithms. These models comprehensively consider the influencing factors of the spatial accessibility of Sanya tourist attractions from different angles, including physical accessibility, tourists' preferences, tourists' choices and spatio-temporal behaviors, so as to reduce repeated calculation and invalid calculation, thus saving calculation time.



Fig. 2. Measurement time of spatial accessibility of tourist attractions in Sanya using different methods.

## V. DISCUSSION

According to the experimental results, the spatial accessibility of Nanwan Monkey Island Ecotourism Area is the best, while the spatial accessibility of Coconut Dream Corridor is the worst. This shows that when planning tourist routes or promoting tourist attractions in Sanya, priority should be given to those scenic spots with good spatial accessibility, such as Nanwan Monkey Island Ecotourism Area and Wuzhizhou Tourist Scenic Area. For scenic spots with poor spatial accessibility, such as Coconut Dream Corridor, more resources need to be put into improvement, such as adding traffic lines and improving the service level of public transportation. The spatial accessibility measurement algorithm of Sanya tourist attractions based on seasonal factor adjustment analysis proposed in this paper has high accuracy and reliability. Therefore, in practical application, this method should be given priority to improve the accuracy and efficiency of spatial accessibility measurement of Sanya tourist attractions. With the increase of the number of time periods, the time for measuring the spatial accessibility of Sanya tourist attractions with different methods increases.

This shows that in real-time monitoring or high frequency measurement, we need to pay more attention to computational efficiency and resource consumption. Therefore, in practical application, the appropriate number of time periods should be selected according to the specific situation to balance the calculation efficiency and measurement accuracy.

## VI. CONCLUSION

In this paper, the spatial accessibility measurement algorithm of Sanya tourist attractions based on seasonal factor adjustment analysis is proposed. Through the study of the spatial accessibility of Sanya tourist attractions, this paper puts forward the method of seasonal factor adjustment analysis. By extracting the data characteristics of Sanya tourist attractions in different seasons, combining with the spatial information of Sanya tourist attractions and GIS technology, this paper analyzes the spatial accessibility of Sanya tourist attractions, constructs the spatial accessibility measurement model of Sanya tourist attractions, and realizes the spatial accessibility measurement of Sanya tourist attractions. This method has a good effect of measuring the spatial accessibility of Sanya tourist attractions, which can effectively eliminate the influence of seasonal factors, improve the accuracy of accessibility measurement and shorten the time of accessibility measurement, so as to measure the spatial accessibility of tourist attractions more objectively and accurately. This not only provides a scientific basis for the development of tourism in Sanya, but also helps to improve the tourist experience.

However, there are still some limitations in this study. For example, there may be some deviation when collecting data, and the seasonal influencing factors of some scenic spots may be complicated. Therefore, future research can further expand the data sources and enhance the accuracy of data, and at the same time, deeply study the seasonal influencing factors of different scenic spots, so as to adjust the analysis methods more finely. In addition, we will explore more efficient and accurate measurement algorithms by combining research results in other fields, such as artificial intelligence and machine learning. At the same time, we can further expand the research scope and combine the development of Sanya tourism with environmental protection, cultural inheritance and other related fields to promote the sustainable development of Sanya tourism. Through continuous in-depth study and improvement of seasonal factor adjustment and analysis methods, we hope to provide more valuable scientific basis and decision support for the prosperity and development of Sanya tourism.

### REFERENCES

[1] Tampubolon, F., & Sinulingga, J. (2021). Socialization of Efforts to Increase Environmental Awareness in Pangambatan Village as A Tourist Attraction in Karo Regency. ABDIMAS TALENTA: Jurnal Pengabdian Kepada Masyarakat, 6(1), 91-98. DOI: 10.32734/ABDIMASTALENTA.V6I1.5395.

[2]    Wang, Y. W., Wu, X. Y., Liu, Z. Z., Chen, H., & Zhao, Y. Y. (2022). Spatial Patterns of Tourist Attractions in the Yangtze River Delta Region. Land, 11(9): 1523. DOI: 10.3390/land11091523.

[3]    Yuliviona, R., Azliyanti, E., Tasri, E. S., & Lindawati. (2021). The effect of tourist attraction, location and promotion toward local tourist decision visit to air manis beach in padang city in new normal policy. IOP Conference Series: Earth and Environmental Science, 747(1), 012085. DOI: 10.1088/1755-1315/747/1/012085.

[4]    Kim, G. S., Chun, J., Kim, Y., & Kim, C. K. (2021). Coastal tourism spatial planning at the regional unit: Identifying coastal tourism hotspots based on social media data. ISPRS International Journal of Geo-Information, 10(3), 167-177. DOI: 10.3390/IJGI10030167.

[5]    Zhou, H. T, Ma, Y. S., Fan, Y. Y., & Ning, X. L. (2023). Spatial distribution and accessibility analysis of red tourism resources in Inner Mongolia. Arid Land Geography, 46(5): 814-822. DOI: 10.12118/ j.issn.1000-6060.2022.423.

[6]    Wang, H., & Yang, L. (2022). Study on Spatial Distribution Characteristics and Accessibility of 4 A Level and Above Scenic Spots in Xinjiang. Journal of Sichuan Normal University(Natural Science),45(6):817-829. DOI: 10.3969/j.issn.1001-8395.2022.06.017.

[7]    Liu, W., Dong, A. R., Deng, L., Zhu, T., & Pi, Y. X. (2021). Research on Fairness Measurement of Urban Public Transportation Accessibility Based on Network Open Data. Journal of Wuhan University of Technology(Transportation Science & Engineering),45(06):1045-1050. DOI: 10.3963/j.issn.2095-3844.2021.06.008.

[8]    Dong, L., Lv, Y., Sun, H., Zhi, D., & Chen, T. (2021). GPS Trajectory-Based Spatio-Temporal Variations of Traffic Accessibility under Public Health Emergency Consideration. Journal of Advanced Transportation,2021,(3):1-22. DOI: 10.1155/2021/8854451.

[9]    Pilario, K. E., Tielemans, A., & Mojica, E. R. E. (2022). Geographical discrimination of propolis using dynamic time warping kernel principal components analysis. Expert Systems with Applications, 187, 115938. DOI: 10.1016/j.eswa.2021.115938.

[10]   Zhao, L., Zhao, X., Zhou, H., Wang, X., & Xing, X. (2021). Prediction model for daily reference crop evapotranspiration based on hybrid algorithm and principal components analysis in Southwest China. Computers and Electronics in Agriculture, 190, 106424. DOI: 10.1016/j.compag.2021.106424.

[11]   Geng, L., Li, H., & Liu, L. (2022). Mining method of fuzzy frequent item set based on principal component analysis. Computer Simulation, 39(02), 410-413. DOI: 10.3969/j.issn.1006-9348.2022.02.078.

[12]   Amaro-Mellado, J. L., Melgar-García, L., Rubio-Escudero, C., & Gutiérrez-Avilés, D. (2021). Generating a seismogenic source zone model for the Pyrenees: A GIS-assisted triclustering approach. Computers & Geosciences, 150, 104736. DOI: 10.1016/ j.cageo.2021.104736.

[13]   Zhou, Z., Qingshan, Z., Dongyi, L., & Weihong, T. (2021). Three-dimensional reconstruction of huizhou landscape combined with multimedia technology and geographic information system. Mobile Information Systems, 2021, 1-13. DOI: 10.1155/2021/9930692.

[14]   Özer Genç, Ç., & Arıcak, B. (2022). Developing a Harvest Plan by Considering the Effects of Skidding Techniques on Forest Soil Using a Hybrid TOPSIS-Entropy Method. Forest Science, 68(3), 312-324. DOI: 10.1093/forsci/fxac010.

[15]   Yan, X. F. (2022). Research on the action mechanism of circular economy development and green finance based on entropy method and big data. Journal of Enterprise Information Management, 35(4), 988-1010. DOI: 10.1108/jeim-01-2021-0024.

# Decoding the Narrative: Patterns and Dynamics in Monkeypox Scholarly Publications

Muhammad Khahfi Zuhanda[1], Desniarti[2], Anil Hakim Syofra[3],
Andre Hasudungan Lubis[4], Prana Ugiana Gio[5], Habib Satria[6], Rahmad Syah[7]

Faculty of Engineering, Universitas Medan Area, Medan, Indonesia[1, 4, 6, 7]
Faculty of Teaching and Educational Sciences, Universitas Muslim Nusantara Al Washliyah, Medan, Indonesia[2]
Faculty of Teaching and Educational Sciences, Universitas Asahan, Kisaran, Indonesia[3]
Department of Mathematics, Universitas Sumatera Utara, Medan, Indonesia[5]

*Abstract*—This study conducts a bibliometric analysis of monkeypox research to uncover trends, influential publishers, and key research topics. A dataset of Google Scholar-indexed articles was analyzed using bibliometric methods and tools such as Publish or Perish (PoP), VOSviewer, and Bibliometrix. The study reveals a growing research interest in monkeypox, with a notable increase in publications over the past decade. The Wiley Online Library emerged as the leading publisher, while highly cited articles covered various aspects of the disease. Cluster analysis identified key research topics, including clinical features, zoonotic transmission, and outbreak patterns. Network visualization and bigram analysis showcased relationships between authors, keywords, and publishers, with "monkeypox" being the most frequent keyword. By visualizing topic trends over time, the study identified emerging areas of investigation. The findings contribute to a comprehensive understanding of monkeypox research, aiding in identifying research gaps and guiding future studies. This research highlights the relevance of bibliometric analysis in health and information sciences. By uncovering trends, influential publishers, and key topics in monkeypox research, this study informs prevention, vaccination, and treatment strategies for mitigating the impact of monkeypox on public health.

*Keywords*—*Bibliometrics; monkeypox virus; research trends; publication patterns; research impact*

## I. INTRODUCTION

Health problems have become a concern lately, still not recovering from COVID-19, which has haunted humans since 2019 [1]. Now in 2022, the outbreak of the monkeypox virus is a scourge that is quite scary for humans. Monkeypox is a disease of animals that mutate to be able to transmit to humans. Monkeypox is caused by the MPXV virus, which is a member of the genus Orthopoxvirus and the family Poxviridae. This case initially spread from the Democratic Republic of Congo to several parts of the world, such as America, Asia, and Europe. Even the death rate caused by the MPXV virus for unvaccinated cases reaches 10% [2].

With the spread of this virus variant in various parts of the world, more research on this monkeypox virus is to find solutions for prevention [3], vaccination [4], and treatment [5]. The health impact is devastating in terms of health and the economy [6]. Preventive measures are needed to solve this problem. Research needs a method to find things that are being studied a lot and map the topic and author to the geographical location of the researcher. Looking for trending topics by geography, time, etc., is essential.

From the literature search, there are still not many studies regarding the bibliometric analysis of monkeypox cases. This research is significant due to the scarcity of studies focusing on the bibliometric analysis of monkeypox cases. The lack of research in this area highlights the importance of conducting this study to reveal trends and insights related to monkeypox over the years.

Despite the urgency of this matter, there exists a scarcity of studies specifically focusing on bibliometric analysis of monkeypox cases. This research aims to fill this critical gap by employing advanced methodologies, including word cluster analysis, network visualization, overlays, and density using VOSviewer software. Additionally, bigram analysis is utilized to examine the relationships between authors, words, and publishers through Bibliometrix. These information science methods, particularly in bibliometrics, play a pivotal role in analyzing patterns, trends, and relationships within scientific literature. The objective of this study is to comprehensively map the scientific landscape, identify research gaps, and provide a profound understanding of the monkeypox domain through the application of these methodologies to monkeypox research.

## II. LITERATURE REVIEW

Bibliometric analysis is a scientific method that integrates mathematical and statistical approaches to analyze and visualize data to determine the structure of topics periodically, develop models, and seek research priorities in a particular field [7]. Bibliometrics research is widely used in the health sector, but according to this paper observation, there have not been many studies involving monkeypox bibliometrics analysis. Publish or Perish (PoP) is software that makes it easy to mine database data from various sources, such as databases indexed by Scopus, Web of Science, Crossref, PubMed, Google Scholar (GS), and others [8]–[10]. In this study mine from the GS-indexed database. GS is a service that makes it easy to collect scientific publications indexed by GS. GS also provides a database that can be used for scientific purposes. Unlike similar services such as Scopus, Web of Science, Pubmed, and others, which require a subscription to access

their databases, GS can be used complementary to mine databases indexed by GS [11].

Many publications use this software [12]–[14]. Arias-Chávez et al. [15] using PoP to mine publication data relating to global scientific production on social networks during the Covid-19 pandemic. Postigo-Zumarán et al. [16] used PoP to search for publications on world scientific production on education and COVID-19 between January 2020 to September 2021. Data from the research of both publications came from five databases, namely Scopus, WoS, GS, Microsoft Academic, and Crossref.

VOSviewer can process data in RIS, JSON, and TXT files. Researchers have widely used VOSviewer to visualize literature data, as has been done by Huang, Z. et al. [17], Kirkendall and Krustrup [18], Pagan-Castaño et al. [19] ,and Papatsounis et al. [20]. Bamel et al. [21] examines leading publication trends over a two-decade period (2000-2020), researching the extent and impact of intellectual capital research in the Journal of Intellectual Capital (JIC). Huang, X. et al. [22] investigates related research advances in pharmaceutical science and pharmaceutical education from a bibliometric point of view and aims to provide advice in facilitating the development of pharmaceutical science and pharmacy graduate education [23].

Bibliometrix is open-source software for comprehensive scientific mapping analysis of scientific literature. Mayara et al. [24] build the map the research literature on Biochemistry education, indexed on the Web of Science covering the scientific production, build upon a recent method to simplify some of the key steps of merging datasets when using the R package Bibliometrix to perform bibliometric analyses. Andrea [25] built a new method to simplify some of the key steps of combining data sets when using the Bibliometrics R package to perform bibliometric analysis.

The significance of this research extends beyond the realm of health sciences and biomedical research. It also highlights the relevance of information sciences in investigating and visualizing trends and patterns in health-related research. By leveraging bibliometric analysis, this study demonstrates the applicability of information science methods in analyzing and synthesizing large volumes of scientific literature. This research provides valuable insights into the field of monkeypox, enabling researchers to gain a deeper understanding of the topic and identify areas that require further investigation.

Furthermore, by employing bibliometric analysis, this study contributes to the broader field of information sciences by showcasing the power of these methodologies in uncovering hidden patterns, revealing research trends, and facilitating evidence-based decision-making. Overall, this research is a significant step in bridging the gap between health sciences, biomedical research, and information sciences, highlighting the interdisciplinary nature of scientific inquiry and its potential for meaningful contributions to public health and knowledge management.

## III. METHODOLOGY

### A. Data Collections and Processing

This study relies on a meticulous and systematic approach to collecting and processing article data indexed by Google Scholar (GS). Unlike subscription-based databases like Scopus or Web of Science, GS provides open access to its data, eliminating the need for researchers to subscribe. This open-source nature enhances accessibility and inclusivity in obtaining data for research purposes.

To ensure the reliability and comprehensiveness of the study, we leverage the PoP (Publish or Perish) feature of Google Scholar. PoP is an open-source application that facilitates the collection of journal data not only from GS but also from various reputable sources such as Scopus, Web of Science, PubMed, Crossref, and Semantic Scholar. The data collected is systematically stored in *.ris format, a widely accepted format for bibliographic citation data, enabling seamless processing in subsequent stages of bibliometric analysis.

For data visualization and in-depth analysis, we employ two powerful tools: VOSviewer and Bibliometrix. VOSviewer is instrumental in creating visual maps based on various parameters such as words, authors, and publishers. This visualization is crucial for identifying emerging research topics, discerning gaps in the literature, and determining the novelty of a study [26]–[29].

Bibliometrix, designed as open-source software integrated with the R programming language, plays a pivotal role in conducting comprehensive scientific mapping analysis. It allows the import of data from various databases, including Scopus, Web of Science, PubMed, Crossref, and others [30]–[32].

### B. Bibliometrics Analysis

Impact of specific articles: To discern the primary contributors in the field, this paper meticulously compiled data on monkeypox-related articles using the Publish or Perish software, which facilitates a nuanced examination of academic publications.

This research focuses on uncovering trends, influential publishers, and key research topics within the field of monkeypox. Cluster analysis was conducted to identify common themes and interconnections among research studies [33]. Network visualization techniques were employed to visualize relationships between authors, keywords, and publishers [34]. The collected data were analyzed and interpreted to reveal patterns, emerging topics, and prominent contributors within the monkeypox research domain. Insights were drawn from the network visualization, cluster analysis, and keyword co-occurrence analysis to understand the research landscape comprehensively.

## IV. RESULTS

### A. Annual Publications

From the results of the collection using PoP, 752 articles were collected in that period. The data is saved in *.ris format and extracted with VOSviewer and Bibliometrix. Bibliometric

analysis using VOSviewer visualizes a network of frequently occurring topics with an occurrence value of 20 from the title and abstract of the article. From the processing results, it can be from the results of collecting data all the time using PoP from GS, 313 articles entitled monkeypox were obtained. Fig. 1 shows the number of publications from year to year. It can be seen that in the last ten years, the number of publications that have a topic entitled monkeypox has increased. In 2022 the number of publications increased sharply to 144 from only six articles in 2021. The reappearance of monkeypox caused an increase in publication at this time. The highest number of publications on monkeypox topics in 2022 was in 2017, with 12 publications. In 2010 and 2015, there were 11 publications using monkeypox titles.



Fig. 1. Distribution of the publications per year.

### B. Publications' Patterns

Table I illustrates the leading publishers contributing to the field of monkeypox research, highlighting the number of articles they have published on this specific topic. At the forefront, Wiley Online Library takes the lead with an impressive 26 articles, closely followed by Academic Oxford with 25 articles. The American Society for Microbiology secures the third position, having published 24 articles. Notable contributions are also observed from the National Center for Biotechnology Information (22 articles) and The American Journal of Tropical Medicine and Hygiene (20 articles). Elsevier, a well-established publisher, has contributed significantly with 19 articles, while The BMJ, Europe PMC, and The International Journal of Infectious Diseases have also made substantial contributions with 16, 15, and 10 articles, respectively. PubMed NCBI, a widely used database, also features in the top 10 publishers with 10 articles. Collectively, this comprehensive overview sheds light on the pivotal role these publishers play in disseminating valuable insights and knowledge within the realm of monkeypox research.

Table II enumerates the top 10 articles with the title "Monkeypox," presenting essential details such as citation counts, authors, and publication years. Topping the list is "The detection of monkeypox in humans in the Western Hemisphere" by Reed K.D. et al. (2004), acknowledged with a

remarkable 655 citations. Following closely is Likos A.M. et al.'s (2005) "A tale of two clades: monkeypox viruses," recognized with 296 citations, emphasizing the diversity within monkeypox viruses. The third position is secured by Hooper J.W. et al.'s (2004) "Smallpox DNA vaccine protects nonhuman primates against lethal monkeypox," renowned for its findings on the protective potential of a smallpox DNA vaccine with 240 citations. "Poxvirus dilemmas—monkeypox, smallpox, and biologic terrorism" by Breman J.G. et al. (1998) claims the fourth spot with 223 citations, delving into the complexities of poxviruses and their implications in biologic terrorism. Lastly, "Human monkeypox: clinical features of 282 patients" authored by Ježek Z. et al. (1987) completes the top five, boasting 175 citations for its insights into the clinical features of human monkeypox patients. Collectively, these top-cited articles constitute pivotal contributions to the understanding of monkeypox, showcasing their enduring impact on scholarly discourse and research endeavors.

TABLE I. THE TOP 10 PUBLISHERS WHO PUBLISHED ARTICLES WITH THE TOPIC MONKEYPOX

| Sources | Publisher | Articles |
|---|---|---|
| 1 | Wiley Online Library | 26 |
| 2 | Academic Oxford | 25 |
| 3 | American Society for Microbiology | 24 |
| 4 | National Center for Biotechnology Information | 22 |
| 5 | The American Journal of Tropical Medicine and Hygiene | 20 |
| 6 | Elsevier | 19 |
| 7 | The BMJ | 16 |
| 8 | Europe PMC | 15 |
| 9 | The International Journal of Infectious Diseases | 10 |
| 10 | PubMed NCBI | 10 |

TABLE II. THE TOP 10 ARTICLES WITH THE TITLE MONKEYPOX

| Num. | Title | Cited | Author | Year |
|---|---|---|---|---|
| 1 | The detection of monkeypox in humans in the Western Hemisphere | 655 | Reed K.D. et al. [35] | 2004 |
| 2 | A tale of two clades: monkeypox viruses | 296 | Likos A.M. et al. [36] | 2005 |
| 3 | Smallpox DNA vaccine protects nonhuman primates against lethal monkeypox | 240 | Hooper J.W. et al. [37] | 2004 |
| 4 | Poxvirus dilemmas—monkeypox, smallpox, and biologic terrorism | 223 | Breman J.G. et al. [38] | 1998 |
| 5 | Human monkeypox: clinical features of 282 patients | 175 | Ježek Z. et al. [39] | 1987 |

### C. Topics Cluster

Table III provides a comprehensive overview of topics clusters in the context of monkeypox research. The clustering of topics in the analysis reveals seven distinct thematic clusters prevalent in the literature on monkeypox. In Cluster 1, the focus is on the analysis of human monkeypox cases and

outbreaks. The cluster encompasses discussions on clinical features, detection, diagnosis, and patient-related aspects. Geographical references to regions like Congo, Democratic Republic, the United States, and the Western Hemisphere are also prominent, emphasizing the global nature of the disease.

TABLE III. CLUSTERS OF TOPICS IN MONKEYPOX RESEARCH

| Cluster | Topics |
|---|---|
| 1 | Case, clinical feature, congo, democratic republic, detection, diagnosis, drc, human monkeypox, human monkeypox case, human monkeypox infection, June, may, monkeypox, monkeypox case, outbreak, patient, patient, person, republic, United States, western hemisphere, year, Zaire. |
| 2 | Cynomolgus macaque, disease, evaluation, family poxciridae, genus orthopoxvirus, human, infection, macaque, monkeypox virus, monkeypox virus challenge, monkeypox virus infection, mpxv, nonhuman primate, protection, smallpox, smallpox vaccine, treatment, western Africa, zoonotic disease. |
| 3 | Comparison, cowpox, differentiation, evidence, mpv, orthopoxvirus, vaccinia, vaccinia virus, varicella zoster virus, variola, variola virus, varv, virus. |
| 4 | Animal, central Africa, concern, covid, emergence, first human case, monkeypox disease, Nigeria, spread, west Africa, world. |
| 5 | abstract, country, europe pmc, laboratory, lesson, monkeypox outbreak, study, symptom |
| 6 | Journal, man, monkeypox infection, mpx, transmission, |
| 7 | Africa, Europe, recent outbreak. |

Cluster 2 centers around the impact of monkeypox virus on nonhuman primates. Topics within this cluster include disease evaluation, infection in macaques, discussions on smallpox vaccine, and the broader context of zoonotic disease. This cluster delves into the complex interplay between the virus and nonhuman primates, providing valuable insights into potential transmission dynamics. The third cluster involves a comparative analysis of orthopoxviruses, drawing parallels between monkeypox and related viruses such as vaccinia, variola, and cowpox. Discussions within this cluster focus on differentiation, evidence, and the distinct characteristics of these viruses. This comparative perspective contributes to a nuanced understanding of the broader orthopoxvirus family.

In Cluster 4, the analysis shifts towards the ecological impact of monkeypox, specifically in animals. Central Africa and Nigeria are focal points, addressing concerns related to the emergence of monkeypox, its potential connection to COVID, and its spread in West Africa and globally. This cluster underscores the broader implications of the disease within the animal population. The fifth cluster is centered on the study of monkeypox outbreaks and the lessons derived from them. Topics within this cluster include the abstracts of studies, lessons learned, symptoms, and country-specific analyses. This cluster provides valuable insights into the dynamics of monkeypox outbreaks and serves as a resource for understanding and managing future occurrences.

Cluster 6 focuses on journal publications related to monkeypox, highlighting aspects such as infection in humans (man), transmission dynamics, and the role of journals in disseminating information. This cluster sheds light on the dissemination of knowledge within the academic community and the dynamics of information exchange. Lastly, Cluster 7

explores the perspectives on monkeypox in Africa and Europe. This geographical focus includes discussions on recent outbreaks, providing insights into how monkeypox is perceived and managed in these regions. The cluster captures the regional nuances and highlights the importance of considering diverse perspectives in understanding and addressing monkeypox.

### D. Relationship Analysis

In Fig. 2 present a detailed analysis of the intricate relationship between authors, abstract content, and publishers using Bibliometrix. The network visualization reveals a rich tapestry consisting of 19 authors, 20 frequently occurring words in the abstracts, and 20 prominent publishers.

Among the 20 words recurrently found in the abstracts, "monkeypox" emerges as the most prevalent, appearing a remarkable 444 times. Other frequently occurring terms include "virus" (22 times), "human" (78 times), "infection" (45 times), "outbreak" (45 times), "disease" (37 times), "smallpox" (29 times), "mpxv" (27 times), "Africa" (25 times), "mpx" (24 times), "variola" (19 times), "zoonotic" (19 times), and "orthopoxvirus" (17 times).

The prevalence of "monkeypox" in abstracts is further underscored by 19 authors consistently incorporating this term in their contributions. Notably, Ii Zuka stands out, utilizing 13 of the identified frequently occurring words a total of 29 times. Following closely are M Saijo, who employs 10 words 25 times, and Ci Hutson, who incorporates 9 words 21 times.

In the realm of publishers, ASM emerges as a leading contributor, publishing articles that prominently feature the identified frequently used words a substantial 151 times. WOL (Wiley Online Library) follows closely with 142 instances, Academic Oxford with 124 instances, NCBI (National Center for Biotechnology Information) with 112 instances, and Elsevier with 89 instances. These findings highlight the pivotal role of these publishers in disseminating research that revolves around the recurrent themes identified in the abstracts, providing key insights into the landscape of monkeypox literature.



Fig. 2. Sankey diagram relationship between author, abstract, and publisher.

## E. Bigram Analysis

To conduct a comprehensive analysis of co-occurring terms within the abstracts of 313 identified articles, this paper employed bigram analysis. Fig. 3 illustrates the noteworthy findings from this analysis. Notably, the phrase "monkeypox virus" emerges as the most frequently occurring bigram, appearing a substantial 117 times, constituting 22% of the total occurrences. Following closely is "human monkeypox," observed in 71 instances, comprising 13% of all identified bigrams.



Fig. 3.    The most frequent bigrams.

The subsequent five most prevalent bigrams include "virus mpxv" (25 occurrences, 5%), "monkeypox outbreak" (23 occurrences, 4%), "monkeypox mpx" (20 occurrences, 4%), "democratic republic" (15 occurrences, 3%), and "zoonotic disease" (14 occurrences, 3%). This insightful analysis sheds light on the prominent word associations within the abstracts, offering valuable insights into prevalent themes and connections in the realm of monkeypox research.

## F. Mapping of Monkeypox Topics

Fig. 4 shows the network of keywords. The red color indicates the topics in cluster 1, the green color represents the members in cluster 2, the blue color represents the terms in cluster 3, and the yellow color represents the topics in cluster 4, the purple color represents the members in cluster 5, the light blue color represents the members in cluster 6, the orange color represents the members in cluster 7. The topic "monkeypox" is the topic that most often appears in cluster 1 with a total number of link strength (TLS) of 874 and occurrence of 214. In cluster 2, "monkeypox virus" became the center of the cluster with a TLS value of 469 and occurrences of 96. The topic "virus" became the center of cluster 3 with a TLS value of 112 and the number of occurrences 24. The topic "animal" is the center of cluster 4 with the number of TLS 65 and the number of occurrences 11. The topic "monkeypox outbreak" is the center of cluster 5 with the number of TLS 106 and the number of occurrences 29. The topic "mpx" is the center of cluster 6 with the number of TLS 136 and the number of occurrences 20. The topic "africa" is the center of cluster 7 with the number of TLS 80 and the number of occurrences 17.



Fig. 4.    Network visualization of monkeypox topics.



Fig. 5.    Overlay visualization of monkeypox topics.



Fig. 6.    Density visualization of monkeypox topics.

Fig. 5 maps the trend of the topic network regarding the "monkeypox virus" from 1964 to 2022. Publications using the title "monkeypox" are mapped from 1964 to 2022. The yellow topics are the most recently discussed with the keyword "monkeypox ". Fig. 6 is a density visualization, and this image shows that the lighter the yellow color and the larger the circle size, the stronger the link with the keyword. So the topic

"monkeypox virus", "monkeypox virus", "case", "infection", and "congo" is a general topics with the keywords "monkeypox virus". Meanwhile, if the color fades and shrinks mixed with green, this is not a common topic to discuss. So the topics "europe," "family poxciridae ", and " recent outbreak" become topics that are not very strongly related to topics.

## V. DISCUSSIONS

The results of bibliometric analysis align with existing literature and contribute valuable insights into the publication patterns and trends in monkeypox research. The upward trajectory in annual publications, particularly the significant surge in 2022, corroborates with the findings of Bunge et al., who noted the changing epidemiology of human monkeypox as a potential threat [2]. The correlation between the rise in publications and the resurgence of monkeypox cases emphasizes the dynamic nature of this emerging public health issue, echoing discussions on the evolving landscape of infectious diseases such as COVID-19 [1].

Examining publishers in the field, this study identifies Wiley Online Library (WOL), Academic Oxford, and the American Society for Microbiology (ASM) as key contributors, aligning with the diverse range of publishers identified in various bibliometric analyses [7]. This information provides researchers with valuable insights into where to find relevant studies on monkeypox, similar to the considerations made by researchers in exploring information science publications [8].

The identification of highly cited articles in this study, such as "The Detection of Monkeypox in Humans in the Western Hemisphere" [35], reinforces the importance of foundational research in shaping the discourse within the field, as highlighted by the significance of top-cited articles in various scientific domains [7]. These influential papers serve as crucial references and contribute to the collective knowledge base on monkeypox.

The exploration of topic clusters aligns with the systematic review conducted by Bunge et al., which underscores the importance of understanding various aspects of monkeypox, including clinical features and zoonotic transmission [2]. These findings further emphasize the interconnectedness between different research topics within the field, similar to the identification of trends, patterns, and collaborations in nursing career research through bibliometric analysis [7].

The network visualization and bigram analysis align with the broader literature on bibliometric analyses, showcasing the significance of specific keywords such as "monkeypox" and the role of authors and publishers in contributing to the body of knowledge [11]. Similar network analyses have been conducted in various fields, from nursing career research to the exploration of pharmaceutical science and education [7], [22].

Finally, the visualization of topic trends over time contributes to the understanding of the evolving research interests related to the monkeypox virus. This mirrors the approach taken in other bibliometric analyses, where visualizing trends aids in identifying emerging areas of investigation [19], [29].

## VI. CONCLUSIONS

In conclusion, this research on the bibliometric analysis of monkeypox provides valuable insights into the trends and patterns of research in the field. The study reveals a significant increase in the number of publications on monkeypox over the last decade, indicating a growing research interest in the disease. The Wiley Online Library, Academic Oxford, and the American Society for Microbiology are identified as the top publishers in this area. The most cited articles on monkeypox cover various aspects of the disease, including detection, virus characterization, vaccination, and clinical features. These highly cited articles serve as foundational research in the field.

The analysis of topic clusters highlights the main themes and areas of focus in monkeypox research, such as clinical features, zoonotic transmission, animal hosts, and outbreak patterns. This information helps researchers identify research gaps and explore interconnected topics within the field. The network visualization and bigram analysis provide insights into the relationships between authors, abstract keywords, and publishers. The word "monkeypox" emerges as the most frequently occurring keyword, emphasizing its central importance in the field. Furthermore, the visualization of topic trends over time demonstrates the evolving research interests related to monkeypox, allowing researchers to identify emerging areas of investigation. Overall, this bibliometric analysis serves as a foundation for future research in monkeypox and provides researchers with valuable information for guiding their studies, identifying research gaps, and contributing to the understanding and managing of the disease.

## REFERENCES

[1] M. K. Zuhanda et al., "Supply chain strategy during the COVID-19 terms: sentiment analysis and knowledge discovery through text mining," Indonesian Journal of Electrical Engineering and Computer Science, vol. 30, no. 2, p. 1120, May 2023, doi: 10.11591/ijeecs.v30.i2.pp1120-1127.

[2] E. M. Bunge et al., "The changing epidemiology of human monkeypox—A potential threat? A systematic review," PLoS Negl Trop Dis, vol. 16, no. 2, 2022, doi: 10.1371/journal.pntd.0010141.

[3] H. Harapan et al., "Physicians' willingness to be vaccinated with a smallpox vaccine to prevent monkeypox viral infection: A cross-sectional study in Indonesia," Clin Epidemiol Glob Health, vol. 8, no. 4, 2020, doi: 10.1016/j.cegh.2020.04.024.

[4] H. Harapan et al., "Acceptance and willingness to pay for a hypothetical vaccine against monkeypox viral infection among frontline physicians: A cross-sectional study in Indonesia," Vaccine, vol. 38, no. 43, 2020, doi: 10.1016/j.vaccine.2020.08.034.

[5] I. B. Abubakar et al., "Traditional medicinal plants used for treating emerging and re-emerging viral diseases in northern Nigeria," Eur J Integr Med, vol. 49, 2022, doi: 10.1016/j.eujim.2021.102094.

[6] N. Berthet et al., "Genomic history of human monkey pox infections in the Central African Republic between 2001 and 2018," Sci Rep, vol. 11, no. 1, 2021, doi: 10.1038/s41598-021-92315-8.

[7] O. Bilik, H. T. Damar, G. Ozdagoglu, A. Ozdagoglu, and M. Damar, "Identifying trends, patterns, and collaborations in nursing career research: A bibliometric snapshot (1980–2017)," Collegian, vol. 27, no. 1, 2020, doi: 10.1016/j.colegn.2019.04.005.

[8] M. Lambovska and D. Todorova, "'Publish and flourish' instead of 'publish or perish': A motivation model for top-quality publications,"

Journal of Language and Education, vol. 7, no. 1, 2021, doi: 10.17323/jle.2021.11522.

[9] H. Costa, F. L. do Canto, and A. L. Pinto, "Google scholar metrics and new the new qualis evaluation system: Impact of the brazilian journals of information science," Informacao e Sociedade, vol. 30, no. 1. 2020. doi: 10.22478/ufpb.1809-4783.2020v30n1.50676.

[10] T. Kwanya, "Publishing and perishing? Publishing patterns of information science academics in Kenya," Information Development, vol. 36, no. 1, 2020, doi: 10.1177/0266666918804586.

[11] H. Soegoto, E. S. Soegoto, S. Luckyardi, and A. A. Rafdhi, "A Bibliometric Analysis of Management Bioenergy Research Using Vosviewer Application," Indonesian Journal of Science and Technology, vol. 7, no. 1, pp. 89–104, 2022, doi: 10.17509/ijost.v7i1.43328.

[12] B. R. Francis, R. bin Ahmad, and S. M. binti Abdullah, "A 59 Years (1962-2021) Bibliometric Analysis of Organizational Support Research Articles," International Journal of Academic Research in Business and Social Sciences, vol. 12, no. 1, 2022, doi: 10.6007/ijarbss/v12-i1/12056.

[13] S. Husin, H. S. Gafur, M. Siscawati, Y. Kristiadi, and A. G. Tangkudung, "Sustainable Palm Oil Industry: Literature Study with Bibliometric Analysis," Budapest International Research and Critics Institute-Journal, vol. 4, no. 4, 2021.

[14] M. T. Picardal and J. M. P. Sanchez, "Effectiveness of Contextualization in Science Instruction to Enhance Science Literacy in the Philippines: A Meta-Analysis," International Journal of Learning, Teaching and Educational Research, vol. 21, no. 1. 2022. doi: 10.26803/ijlter.21.1.9.

[15] D. Arias-Chávez, R. W. A. Espejo, R. C. Juro, J. C. Tovar, and J. E. P. Latour, "Scientific Production on Social Networks during the COVID 19 Pandemic," Webology, vol. 19, no. 1, 2022, doi: 10.14704/web/v19i1/web19144.

[16] J. E. Postigo-Zumarán, C. E. G. Carrión, R. A. S. Alviar, M. B. Q. Calderón, and D. Arias-Chávez, "World Scientific Production on Education and COVID 19: A Bibliometric Analysis," Webology, vol. 19, no. 1, 2022, doi: 10.14704/web/v19i1/web19143.

[17] Z. Huang et al., "Progress on Pharmaceutical Sciences/Pharmacy Postgraduate Education: a Bibliometric Perspective," J Pharm Innov, 2022, doi: 10.1007/s12247-021-09611-z.

[18] D. T. Kirkendall and P. Krustrup, "Studying professional and recreational female footballers: A bibliometric exercise," Scand J Med Sci Sports, vol. 32, no. S1, 2022, doi: 10.1111/sms.14019.

[19] E. Pagan-Castaño, J. C. Ballester-Miquel, J. Sánchez-García, and M. Guijarro-García, "What's next in talent management?," J Bus Res, vol. 141, 2022, doi: 10.1016/j.jbusres.2021.11.052.

[20] A. G. Papatsounis, P. N. Botsaris, and S. Katsavounis, "Thermal/Cooling Energy on Local Energy Communities: A Critical Review," Energies, vol. 15, no. 3. 2022. doi: 10.3390/en15031117.

[21] U. Bamel, V. Pereira, M. del Giudice, and Y. Temouri, "The extent and impact of intellectual capital research: a two decade analysis," Journal of Intellectual Capital, vol. 23, no. 2. 2022. doi: 10.1108/JIC-05-2020-0142.

[22] X. Huang et al., "A Bibliometric Analysis Based on Web of Science: Current Perspectives and Potential Trends of SMAD7 in Oncology," Front Cell Dev Biol, vol. 9, 2022, doi: 10.3389/fcell.2021.712732.

[23] C. Priovashini and B. Mallick, "A bibliometric review on the drivers of environmental migration," Ambio, vol. 51, no. 1. 2022. doi: 10.1007/s13280-021-01543-9.

[24] M. L. de O. Barbosa and E. Galembeck, "Mapping research on biochemistry education: A bibliometric analysis," Biochemistry and Molecular Biology Education, vol. 50, no. 2, 2022, doi: 10.1002/bmb.21607.

[25] A. Caputo and M. Kargina, "A user-friendly method to merge Scopus and Web of Science data during bibliometric analysis," Journal of Marketing Analytics, vol. 10, no. 1, 2022, doi: 10.1057/s41270-021-00142-7.

[26] S. H. H. Shah, S. Lei, M. Ali, D. Doronin, and S. T. Hussain, "Prosumption: bibliometric analysis using HistCite and VOSviewer," Kybernetes, vol. 49, no. 3, 2020, doi: 10.1108/K-12-2018-0696.

[27] I. Hamidah, Sriyono, and M. N. Hudha, "A bibliometric analysis of COVID-19 research using vosviewer," Indonesian Journal of Science and Technology, vol. 5, no. 2, 2020, doi: 10.17509/ijost.v5i2.24522.

[28] X. Ding and Z. Yang, "Knowledge mapping of platform research: a visual analysis using VOSviewer and CiteSpace," Electronic Commerce Research, vol. 22, no. 3, 2022, doi: 10.1007/s10660-020-09410-7.

[29] J. K. Tamala, E. I. Maramag, K. A. Simeon, and J. J. Ignacio, "A bibliometric analysis of sustainable oil and gas production research using VOSviewer," Clean Eng Technol, vol. 7, 2022, doi: 10.1016/j.clet.2022.100437.

[30] K. K. Ingale and R. A. Paluri, "Financial literacy and financial behaviour: a bibliometric analysis," Review of Behavioral Finance, vol. 14, no. 1. 2022. doi: 10.1108/RBF-06-2020-0141.

[31] J. Wang, X. Li, P. Wang, Q. Liu, Z. Deng, and J. Wang, "Research trend of the unified theory of acceptance and use of technology theory: A bibliometric analysis," Sustainability (Switzerland), vol. 14, no. 1, 2022, doi: 10.3390/su14010010.

[32] Z. Xiao, Y. Qin, Z. Xu, J. Antucheviciene, and E. K. Zavadskas, "The Journal Buildings: A Bibliometric Analysis (2011–2021)," Buildings, vol. 12, no. 1, 2022, doi: 10.3390/buildings12010037.

[33] E. Elihami and M. Melbourne, "The Trend of 'Multicultural Education' in 2021-2022: Bibliometrics Mapping in Scopus," Jurnal Pendidikan Progresif, vol. 12, no. 1, 2022, doi: 10.23960/jpp.v12.i1.202104.

[34] X. Li, K. Wu, and Y. Liang, "A Review of Agricultural Land Functions: Analysis and Visualization Based on Bibliometrics," Land (Basel), vol. 12, no. 3, 2023, doi: 10.3390/land12030561.

[35] K. D. Reed et al., "The Detection of Monkeypox in Humans in the Western Hemisphere," New England Journal of Medicine, vol. 350, no. 4, 2004, doi: 10.1056/nejmoa032299.

[36] A. M. Likos et al., "A tale of two clades: Monkeypox viruses," Journal of General Virology, vol. 86, no. 10, 2005, doi: 10.1099/vir.0.81215-0.

[37] J. W. Hooper et al., "Smallpox DNA Vaccine Protects Nonhuman Primates against Lethal Monkeypox," J Virol, vol. 78, no. 9, 2004, doi: 10.1128/jvi.78.9.4433-4443.2004.

[38] J. G. Breman and D. A. Henderson, "Poxvirus Dilemmas — Monkeypox, Smallpox, and Biologic Terrorism," New England Journal of Medicine, vol. 339, no. 8, 1998, doi: 10.1056/nejm199808203390811.

[39] Z. Ježek, M. Szczeniowski, K. M. Paluku, and M. Mutombo, "Human Monkeypox: Clinical Features of 282 Patients," Journal of Infectious Diseases, vol. 156, no. 2, 1987, doi: 10.1093/infdis/156.2.293.

# Healthcare Intrusion Detection using Hybrid Correlation-based Feature Selection-Bat Optimization Algorithm with Convolutional Neural Network

A Hybrid Correlation-based Feature Selection for Intrusion Detection Systems

H. Kanakadurga Bella[1], S. Vasundra[2]

Department of CSE, Jawaharlal Nehru Technological University Anantapur, Ananthapuramu, India[1]
Department of CSE, JNTUA College of Engineering, Ananthapuramu,
Constituent College of Jawaharlal Nehru Technological University Anantapur, Ananthapuramu, India[2]

*Abstract*—Cloud computing is popular among users in various areas such as healthcare, banking, and education due to its low-cost services alongside increased reliability and efficiency. But, security is a significant problem in cloud-based systems due to the cloud services being accessed via the Internet by a variety of users. Therefore, the patient's health information needs to be kept confidential, secure, and accurate. Moreover, any change in actual patient data potentially results in errors during the diagnosis and treatment. In this research, the hybrid Correlation-based Feature Selection-Bat Optimization Algorithm (HCFS-BOA) based on the Convolutional Neural Network (CNN) model is proposed for intrusion detection to secure the entire network in the healthcare system. Initially, the data is obtained from the CIC-IDS2017, NSL-KDD datasets, after which min-max normalization is performed to normalize the acquired data. HCFS-BOA is employed in feature selection to examine the appropriate features that not only have significant correlations with the target variable, but also contribute to the optimal performance of intrusion detection in the healthcare system. Finally, CNN classification is performed to identify and classify intrusion detection accurately and effectively in the healthcare system. The existing methods namely, SafetyMed, Hybrid Intrusion Detection System (HIDS), and Blockchain-orchestrated Deep learning method for Secure Data Transmission in IoT-enabled healthcare systems (BDSDT) are employed to evaluate the efficacy of HCFS-BOA-based CNN. The proposed HCFS-BOA-based CNN achieves a better accuracy of 99.45% when compared with the existing methods: SafetyMed, HIDS, and BDSDT.

*Keywords*—*Convolutional neural network; deep learning; intrusion detection system; healthcare; security*

## I. INTRODUCTION

Network Intrusion Detection Systems (NIDSs) identify malicious activities and safeguard the vulnerable services by monitoring network traffic and providing alerts when anomalous events are recognized. Some organizations that are primarily focused on obtaining private user data, establishing the foundation for modern-day detection and protection are attacked by cyber-attackers. Furthermore, the healthcare sector keeps growing, and most hospitals are integrating e-healthcare systems as quickly as feasible to fulfill the needs of their patients. IDS based on cloud networks employ anomaly-based techniques to protect the cloud-based applications [1]. In network security, there are two common detection techniques for NIDS, anomaly-based detection, and signature-based detection [2]. An anomaly-based IDS analyzes the network traffic and correlates it to a created baseline for unknown or known attacks, where a signature-based IDS is allowed to be employed while the attack patterns are established and pre-determined [3, 4]. To address numerous security issues, the cloud utilizes numerous cybersecurity techniques like IDS, Intrusion Prevention Systems (IPS), and firewalls [5]. The centralized processing technique used by cloud computing involves uploading every transaction and processing the end-user service requests based on the transmission bandwidth, capacity of storage, and computer resources [6]. Proactive network security defenses are required to protect essential assets and data because the cloud attack vector has the potential to result in successful security breaches [7].

Network security has always placed a high priority on intrusion detection since it is crucial for identifying anomalous activity on secured internal networks [8]. The network of intermediate, source, and endpoint are used to identify the Distributed Denial-of-Services (DDoS) attacks. The attack's endpoint is easily detected because of the massive volume of network traffic that is generated [9]. A significant number of traditional intrusion detection systems use either a port-based or Deep Packet Inspection (DPI) technique. The port-based technique identifies traffic by using the ports established by the Internet Assigned Numbers Authority (IANA) [10]. Software Defined Network (SDN) is an emerging design that is cost-effective, flexible, adaptable, and controlled, thereby making it more suitable for presently employed complicated applications and bandwidth [11, 12]. SDN's goal is to create a logically centralized hub for internet and networking architects so that they quickly respond to the evolving client demands [13]. Deep learning techniques, especially CNN represent remarkable capacity in automatically extracting features and intricate patterns from complex data, including network traffic [14]. By employing Deep CNN, the IDS efficiently recognizes anomalous behavior and emerging threats in real time [15]. Since cloud services are accessed via the internet by a variety of users, security is a significant concern in cloud-based systems because the health information of patients must be

kept confidential, secure, and accurate. Moreover, any change in actual patient data results in errors during the diagnosis and treatment [31-34]. Therefore, the HCFS-BOA based on CNN is proposed in this research, for intrusion detection to secure the entire network in the healthcare system. The main contributions of this research are as follows:

- The proposed HCFS-BOA approach is evaluated on the CIC-IDS2017, NSL-KDD benchmark datasets, and the Min-max normalization technique is employed to normalize the raw data.

- For feature selection, HCFS-BOA is employed to examine the appropriate features that not only have significant correlations with the target variable, but also contribute to the optimal performance of intrusion detection in the healthcare system.

- Finally, CNN is employed for classification to identify and classify intrusion detection accurately and effectively. The efficacy of HCFS-BOA is analyzed based on the performance measures of accuracy, precision, recall, and f1-score.

The rest of the paper is organized as follows: Section II presents the literature survey. The block diagram of the proposed method is discussed in Section III. The results are illustrated in Section IV, while Section V discusses the conclusion of this paper.

## II. LITERATURE SURVEY

Faruqui et al. [16] presented a SafetyMed for Internet of Medical Things (IoMT) IDS by employing hybrid CNN-Long Short-Term Memory (CNN-LSTM). The SafetyMed was the first IDS that included an optimization approach based on the trade-off between Detection Rate (DR) and False Positive Rate (FPR). The SafetyMed enhanced the safety and security of medical devices and patient information. However, the presented SafetyMed method had no defense mechanism against an attack of Adversarial Machine Learning (AML).

Vashishtha et al. [17] implemented a HIDM for cloud-based healthcare systems to detect all kinds of attacks. The hybrid approach was a mixture of a Signature-based Detection Model (SDM) and an Anomaly-based Detection Model (ADM). The datasets of NSL-KDD, CICIDS2017, and UNSW-NB15 were employed to evaluate the efficacy of the HIDM approach. The implemented method had a higher detection rate with the error of Type-I and Type-II for both ADM and SDM. However, combining various detection systems increased the risk of false negatives and false positives.

Kumar et al. [18] introduced a BDSDT for the transmission of secure data in IoT-based healthcare systems. Initially, the architecture of blockchain was created in all IoT devices that were identified and established using a zero-knowledge proof, and then connected to the blockchain network using a smart contract-based ePOW consensus. Then, a bidirectional LSTM was employed using a DL to recognize IDS in the healthcare system. The BDSDT enhanced the privacy and security by combining both DL and blockchain methods. However, BDSDT wasn't effective against web and Bot threat attacks as

there were fewer instances of these two classes which led to changes in actual patient data resulting in errors during the diagnosis and treatment.

Halbouni et al. [19] presented a CNN-Long Short-Term Memory (CNN-LSTM) for IDS system. The ability of CNN to extract the spatial features alongside the ability of LSTM to extract the temporal features were the highlights of this model. In order to improve performance, batch normalization and the layers of dropout were created to the presented method. The presented method decreased the false alarm rate and improved the rate of detection. However, CNN-LSTM failed to provide a high detection rate for specific kinds of attacks like web attacks and worms which led to changes in actual patient data resulting in errors during the diagnosis and treatment.

Han et al. [20] presented an Intrusion Detection Hyperparameter Control System (IDHCS) to regulate and train a Deep Neural Network (DNN) extracted feature and the module of k-means clustering in terms of Proximal Policy Optimization (PPO). The most valuable network features were extracted by the DNN under the control of an IDHCS, which also used K-means clustering to detect intrusion. The IDHCS performed effectively for each dataset, as well as the combined dataset. However, to represent a more realistic network environment, a diverse dataset needed to be examined.

Bakro et al. [21] introduced a hybrid feature selection approach that combined filter techniques such as Particle Swarm Optimization (PSO), Chi-Square (CS), and Information Gain (IG). Combining each of these three techniques was a novel method that generated a more reliable process of feature selection by using every technique's strength to increase the possibilities of selecting the most associated features. The introduced method had the benefits of flexibility, time complexity, interpretability, and scalability. But, the feature selection approach was not done properly which resulted in overfitting.

Sudar et al. [22] implemented a Machine Learning (ML) approach based on Decision Tree (DT) and Support Vector Machine (SVM) to detect Distributed Denial of Service (DDoS) attacks. The classification approach was established in the environment of Software Defined Network (SDN). The DT and SVM approaches were deployed to distinguish among malicious and normal traffic data. This approach provided better accuracy and detection rate. Nonetheless, this implemented approach struggled to adapt to evolving attack strategies.

Praveena et al. [23] developed a Deep Reinforcement Learning approach that was optimized by Black Widow Optimization (DRL-BWO) for intrusion detection in Unmannered Aerial Vehicles (UAV). The BWO approach was deployed for parameter optimization of the DRL method which assisted in enhancing the performance of intrusion detection in UAV networks. This approach was fit for the tasks of information extraction in high dimensional space. Nonetheless, the intricate nature of the DRL-BWO approach resulted in minimized interpretability.

Chinnasamy et al. [24] presented a Blockchain DDoS flooding attack with dynamic path detectors. The ML approach

was established to identify the attacks which focused on the DDoS assault. The primary essential traits were employed to predict the accurate DDoS attacks by utilizing a different attribute selection approach. Nevertheless, this presented approach led to severe network congestion which hindered the processing of transactions and slowed down the overall system's performance.

Chinnasamy et al. [25] developed an ML approach for effective phishing attack detection. Based on the input features such as Uniform Resource Locator (URL) and Web Traffic, the link was classified as phishing or non-phishing. This approach was determined by retrieving datasets from ML and phishing cases by employing SVM, Random Forest (RF), and Genetic. Nevertheless, ML approaches in phishing detection struggled to maintain pace with constantly evolving phishing tactics which led to potential delays in identifying the new attacks.

Anupriya et al. [26] implemented an ML approach for fraud account detection. To compute buddy similarity criteria, the adjacency network matrix graph was employed and then new features were acquired by utilizing the Principle Component Analysis (PCA). This was employed to equalize the data and transform it into the classifier in the next phase of cross-validation for training and testing the classifier. Nevertheless, due to imbalanced datasets, this approach struggled with evolving the fraud pattern and generated false positives or negatives.

There are some limitations with the existing methods that are mentioned above such as the methods not being effective in detecting attacks which led to changes in actual patient data resulting in errors during the diagnosis and treatment. In order to overcome these issues, the HCFS-BOA-based CNN is proposed for intrusion detection to secure the entire network in the healthcare system.

## III. PROPOSED METHODOLOGY

In this research, a hybrid CFS-BOA-based CNN approach is proposed for intrusion detection in healthcare systems using deep learning. It includes datasets, min-max normalization, feature selection using HCFS-BOA, classification using CNN, and performance evaluation. The overview of the proposed method is represented in Fig. 1.

### A. Datasets

The proposed HCFS-BOA approach is evaluated on CIC-IDS2017 [27] and NSL-KDD benchmark datasets. The CIC-IDS-2017 dataset includes malicious and normal traffic data that is considered new and does not include an enormous amount of redundant data. It includes eleven new attacks namely, PortScan, Brute Force, DoS, web attacks like SSH, Patator, FTP-Patator, SQL injection, and XSS. It is created by the Canadian Institute for Cybersecurity in 2017, and its 80 features are employed to monitor malicious and benign traffic. The NSL-KDD is an extension of the KDD cup 99 database and contains 41-dimensional vectors with numerical and categorical values. The intrusion attacks in the NSL-KDD database are probe attacks, Remote to a user (R2L), Denial of Service (DoS), and the User to Root attack (U2R). NSL-KDD is an IoT dataset used for model training purposes in healthcare applications.

### B. Pre-processing

After data collection, the normalizing process is established by rescaling the attributes with a uniform contribution. Typically, the data normalization technique addresses two key problems: the presence of outliers and the presence of dominant features. The various methods for normalizing data based on the measures of statistics are examined. Consider the data with $z$ records and $N$ instances, as expressed numerically in Eq. (1).

$$Data = \{p_{i,n}, q_n | i \in z \text{ and } n \in N\} \tag{1}$$

where, $q$ indicates the label of class and $p$ represents the data to be learned via a learning process. The Min-max normalization technique [28] is employed to normalize the raw data, which is one of the various normalization techniques. This approach greatly minimizes the outlier's impact on the data. It scales the obtained data within the range of 0 to 1 which is numerically expressed in Eq. (2).

$$p_{i,n} = \frac{p_{i,n} - min(p_i)}{max(p_i) - min(p_i)}(nMax - nMin) + nMin \tag{2}$$

where, $max$ and $min$ represent the $ith$ attribute's maximum and minimum values. By employing $nMax$ and $nMin$, the acquired data are rescaled by the upper and lower boundaries. This acquired data is then passed as input to the feature selection.



Fig. 1. Block diagram for the proposed method

## C. Feature Selection

After normalizing the acquired data, the hybrid CFS-BOA approach is implemented for feature selection. In CFS-BOA, the features are selected by using a nature-inspired optimization technique to enhance the optimization process. The CFS-BOA's goal is to choose the most useful feature subset for detecting and avoiding security vulnerabilities while minimizing the redundancy and computational complexity. When compared to other optimization algorithms like Ant Colony Optimization (ACO) and Particle Swarm Optimization (PSO), the BOA tunes the optimization process for maximum efficiency for combining with CFS. The HCFS-BOA examines appropriate features that not only have significant correlations with the target variable, but also contribute to its optimal performance of intrusion detection in the healthcare system. This hybrid method has the potential to result in a more efficient and effective IDS that is specific to the unique characteristics of healthcare data and security requirements.

*1) Correlation-based Feature Selection (CFS):* One of the most known filter algorithms is CFS which selects features based on the output of a heuristic (correlation-based) evaluation function. It seeks to choose subsets whose attributes are highly correlated with the class but unassociated with one another. Repetitive features are selected based on their high correlation with at least one other feature, while low-association features are ignored. The function of the CFS feature subset assessment is mathematically expressed in Eq. (3).

$$M_s = \frac{k\overline{r_{cf}}}{\sqrt{k+k(k-1)+\overline{r_{ff}}}} \tag{3}$$

$M_s$ – feature subset's heuristic evaluation for a feature set that includes $k$ features

$\overline{r_{cf}}$ – average degree of connection between the category label and the features

$\overline{r_{ff}}$ – average degree of inter-connection between features

A correlation technique based on the feature subsets is used for the evaluation of CFS. During the procedure, the feature set with the greatest value is determined to decrease the training and testing set size. A larger $\overline{r_{cf}}$ or smaller $\overline{r_{ff}}$ out of the obtained subsets by the approach provides a greater evaluation value.

*2) Bat Optimization Algorithm (BOA):* BOA is the first algorithm for optimization and computational intelligence, influenced by microbat echolocation behavior. In a d-dimensional search, every bat flies at random with $v_i^t$ velocity, $x_i^t$ location and $f_i$ frequency at $t$ iteration. The current best solution $x_*$ is archived for $n$ bats in a population through an iterative search process.

The procedures for updating the location $x_i^t$ and velocity $v_i^t$ at each time step $t$ are mathematically presented in Eq. (4), Eq. (5) and Eq. (6).

$$f_i = f_{min} + (f_{max} - f_{min})\beta \tag{4}$$

$$v_i^t = v_i^{t-1} + (x_i^{t-1} - x_*)f_i \tag{5}$$

$$x_i^t = x_i^{t-1} + v_i^t \tag{6}$$

where, $\beta \in [0,1]$ is a vector selected at random from a uniform distribution.

Once a solution is chosen from the existing ideal solutions, a new solution for every bat is produced via a local random walk which is numerically expressed in Eq. (7).

$$x_{new} = x_{old} + \varepsilon A' \tag{7}$$

where, $\varepsilon$ is a random vector generated from uniform or Gaussian distribution in the range [-1,1].

$A'$ is the average loudness of all bats at a time step.

Furthermore, the rate of pulse emission and loudness are modified as the iterations progress. They are updated using the following Eq. (8) and Eq. (9).

$$A_i^{t+1} = \alpha A_i^t \tag{8}$$

$$r_i^{t+1} = r_i^0(1 - e^{-\gamma t} \tag{9}$$

where, $0 < \alpha < 1$ and $\gamma < 0$ are constant.

*3) HCFS-BOA method for feature selection:* The significance and correlation of the chosen feature subset are evaluated using the HCFS-BOA-based feature selection method. Correlation-based feature method is used in the HCFS-BOA to create a fitness function and assess the reliability of the reduced feature subset. CFS evaluates the correlation of mean feature class and the average inter-correlation between features for $S$ feature subset with $k$ features, where $S = (s_1, s_2, \dots s_k)$ using (3). CFS is a classical filter method that selects relevant features based on correlation-based evaluation due to feature redundancy. By storing solutions in a bat's vector, BA is inspired by the echolocation activity of microbats, eliminating redundant features and reducing dimensionality. When a bat moves, it archives the best solution at the time. During the process of iterative search, the population scans for the optimum arrangement by updating and refreshing the position of each bat based on Eq. (4), Eq. (5), and Eq. (6). An ideal intrusion-detection approach has a higher detection rate and a lower false positive rate. Hence, a weighted fitness function is shown in Eq. (10).

$$Maximize\ Fit = w_1 \times DR \times w_2 \times (1 - FPR) \tag{10}$$

where, $w_1$ and $w_2$ are the weights for the Detection Rate and False Positive Rate, respectively. A higher fitness $Fit$ means higher intrusion detection performance. In one iteration of the HCFS-BOA, the algorithm chooses a feature subset that depends on its correlation coefficients with the target variable. The bat optimization process involves updating the virtual bat's positions in the search space with each bat representing a potential feature subset. The technique iteratively refines feature selection by adjusting the position of bats and evaluates their performances via correlation-based metrics during both the testing and training phases. Thus, the rescaling acquired

data is passed into the feature selection phase which is sufficient for the classification of intrusion detection.

*D. Classification*

The selected features are classified using the CNN model which produces enormous results in domains such as Natural Language Processing (NLP), image processing, and healthcare diagnosis systems. For recognizing patterns and anomalies in network traffic or system logs, CNN classification is employed to improve intrusion detection in healthcare systems. Using CNN classification for IDS in healthcare helps to protect sensitive patient data, ensure the integrity of healthcare information systems, and avoid security breaches. It is an essential component of healthcare cybersecurity measures to protect electronic health records and vital healthcare infrastructure.

In contrast to Multi-Layer Perceptron (MLP), CNN reduces the number of neurons and parameters, resulting in rapid adaptability and minimal complexity. The CNN model offers an extensive number of clinical classification applications. CNN models are a subset of Feed-Forward Neural Network (FFNN) [29, 30] and Deep Learning models. The convolution operations convention is constant which implies that the filter is independent in function, thereby reducing the amount of parameters. Pooling, convolution, and fully connected layers are the three types of layers used in the CNN method. These layers are required for performing feature extraction, dimensionality reduction, and classification. The filter is slid on the computers through the forward pass of convolution operation, and the input capacity of the activation map that assesses the point-wise result of every score is added to obtain the activation. The sliding filter is employed by linear and convolution operators, being stated as a quick distribution of dot product. Consider $w$ is the kernel function, $x$ is the input, $(x \times w)(a)$ at time $t$ is formulated as in Eq. (11).

$$(x \times w)(a) = \int x(t)w(a - t)da \qquad (11)$$

Where, $a$ is $R^n$ for each $n \geq 1$. The parameter $t$ is the discrete which is presented in Eq. (12).

$$(x \times w)(a) = \sum_a x(t)w(t - a) \qquad (12)$$

The 2D image $I$ is given as input, $K$ is a 2D kernel, and the convolution is formulated as in Eq. (13).

$$(I \times K)(i,j) = \sum_m \sum_n I(m,n)K(i - m, j - n) \qquad (13)$$

In order to improve the non-linearity, two activation functions, ReLU and softmax are utilized. The ReLU is mathematically represented as in Eq. (14).

$$ReLU(x) = max(0, x)\ x \in R \qquad (14)$$

The gradient $ReLU(x) = 1$ for $x > 0$ and $ReLU - (x) = 0$ for $x < 0$. The ReLU convergence ability is better than the sigmoid non-linearities. The next layer is softmax, preferable when the result requires including two or more classes which is mathematically formulated as in Eq. (15).

$$softmax(x_i) = \frac{exp(x_i)}{\sum_j exp(x_i)} \qquad (15)$$

The pooling layers are applied to the result in a statistic of input, and the structure of output is rescaled without losing the essential information. There are various types of pooling layers, this paper utilizes the highest pooling that individually produces large values in the rectangular neighborhood of individual points $(i, j)$ in 2D information for every input feature correspondingly. The fully connected (FC) layer, which is the last layer with $m$ and $n$ output and input are illustrated further. The parameter of the output layer is stated as a weight matrix $\in M_{m,n}$. Where, $m$ and $n$ are rows and columns, and the bias vector $b \in R^m$. Consider the input vector $x \in R^n$, the fully connected layer output with an activation function $f$ is formulated as in Eq. (16).

$$FC(x) := f(Wx = b) \in R^m \qquad (16)$$

where, $Wx$ is the matrix product where function $f$ is employed as a component. This fully connected layer is applied for classification difficulties. The FC layer of CNN is commonly involved at the topmost level. The CNN production is compressed and displayed as a single vector.

Table I shows the notation description.

TABLE I. NOTATION DESCRIPTION

| Symbol | Description |
|---|---|
| $q$ | label of class |
| $p$ | data to be learned via a learning process |
| $min$ and $max$ | $ith$ attribute minimum and maximum values |
| $M_s$ | feature subset's heuristic evaluation for a feature set that includes $k$ features |
| $\overline{r_{cf}}$ | the average degree of connection between the category label and the features |
| $\overline{r_{ff}}$ | average degree of inter-connection between features |
| $v_i^t$ | velocity |
| $x_i^t$ | location |
| $f_i$ | frequency at $t$ iteration |
| $\varepsilon$ | random vector generated from a uniform or Gaussian distribution in the range [-1,1] |
| $A'$ | average loudness of all bats at time step |
| $w_1$ and $w_2$ | weights for the Detection Rate and False Positive Rate |
| $Wx$ | Matrix product |

## IV. EXPERIMENTAL RESULTS

In this research, the HCFS-BOA based CNN is simulated using a Python environment with the system configuration of 16GB RAM, Intel core i7 processor, and Windows 10 operating system. The parameters like accuracy, precision, recall, and f1-score are utilized to estimate the performance of the model. The mathematical representation of these parameters is shown Eq. (17) to Eq. (20).

- Accuracy – Accuracy is the proportion of accurate predictions to all input samples and it is calculated using Eq. (10).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \qquad (17)$$

- Precision - The precision measures the percentage of actual data records versus expected data records. The performance of the classification model is greater if the precision is higher.

$$Precision = \frac{TP}{TP+FP} \qquad (18)$$

- Recall – Recall is calculated as the sum of the true positives and the positive class images.

$$Recall = \frac{TP}{TP+FN} \qquad (19)$$

- F1-Score – It is also known as the harmonic mean which seeks a balance between recall and precision.

$$F1 - Score = \frac{2TP}{2TP+FP+FN} \qquad (20)$$

*A. Quantitative and Qualitative Analysis*

This section shows the quantitative and qualitative analysis of the proposed CSF-BOA-based CNN model in terms of precision, accuracy, f1-score, and recall, as presented in Tables II, III and IV. Table II illustrates the performance of feature selection on the CIC-IDS2017 dataset. The performances of ACO, PSO, CFS, and BOA are measured and matched with the proposed HCFS-BOA. Fig. 2 represents a graphical illustration of the feature selection methods. The obtained result shows that the proposed HCFS-BOA algorithm attains an accuracy of 95.98%, precision of 94.23%, recall of 93.62%, and f1-score of 94.96% which is better when compared to the existing optimization algorithms.

Table III illustrates the performance of classification with default features using CIC-IDS2017 dataset. The performance of Support Vector Machine (SVM), Artificial Neural Network (ANN), K-Nearest Neighbor (KNN), and Recurrent Neural Network (RNN) are measured and matched with the proposed HCFS-BOA. Fig. 3 represents the graphical illustration of classification performances. The obtained result shows that the proposed HCFS-BOA algorithm attains an accuracy of 93.68%, precision of 92.92%, recall of 91.69%, and f1-score of 92.73% which is superior when compared to the existing optimization algorithms.

TABLE II. PERFORMANCE OF FEATURE SELECTION USING CIC-IDS2017 DATASET

| Methods | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| ACO | 89.45 | 85.61 | 90.12 | 89.23 |
| PSO | 91.82 | 81.20 | 88.65 | 90.14 |
| CFS | 94.37 | 90.47 | 91.52 | 92.74 |
| BOA | 93.26 | 92.82 | 92.76 | 93.85 |
| HCFS-BOA | 95.98 | 94.23 | 93.62 | 94.96 |



Fig. 2. Graphical representation of feature selection performances.

TABLE III. PERFORMANCE OF CLASSIFICATION WITH DEFAULT FEATURES USING CIC-IDS2017 DATASET

| Methods | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| SVM | 88.21 | 89.65 | 89.33 | 88.37 |
| ANN | 86.45 | 85.37 | 90.27 | 89.91 |
| KNN | 89.98 | 88.83 | 89.24 | 90.28 |
| RNN | 92.47 | 90.61 | 90.49 | 91.96 |
| CNN | 93.68 | 92.92 | 91.69 | 92.73 |

Fig. 3. Graphical representation of classification performances.

Table IV illustrates the classification outcomes with optimized features using CIC-IDS2017 dataset. The performance of SVM, ANN, KNN, and RNN are measured and matched with the optimized feature CNN. Fig. 4 illustrates the graphical representation of classification performances with optimized features. The obtained outcomes prove that the CNN algorithm accomplishes an accuracy of 99.45%, precision of 98.89%, recall of 98.67%, and f1-score of 97.98%, therefore being superior in contrast to the existing optimization algorithms. The ACO, PSO, CFS, and BOA consume 25 seconds, 29 seconds, 31 seconds, and 35 seconds of time, respectively. The time analysis of HCFS-BOA with CNN demands a training time of 20 seconds, being more robust in comparison with other optimization techniques like ACO, PSO, CFS, and BOA on the CIC-IDS2017 dataset. Table V shows the performance of classification with optimized features on the NSL-KDD dataset. Fig. 5 shows that the obtained outcomes of optimized results of the CNN algorithm accomplishes an accuracy of 98.13%, precision of 97.36%, recall of 97.07%, and f1-score of 95.34%, in that way, proving more robust in contrast to the previous optimization algorithms. The ACO, PSO, CFS, and BOA require 22

seconds, 25 seconds, 28 seconds, and 34 seconds of time, respectively. The time analysis of HCFS-BOA with CNN needs a training time of 15 seconds which is lesser than that of the previous optimization techniques like ACO, PSO, CFS, and BOA on the NSL-KDD dataset.



Fig. 4. Graphical representation of optimized feature performances using CIC-IDS 2017.



Fig. 5. Graphical representation of optimized feature performances using NSL-KDD.

TABLE IV. PERFORMANCE OF CLASSIFICATION WITH OPTIMIZED FEATURES USING CIC-IDS2017 DATASET

| Methods | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| SVM | 93.54 | 94.21 | 92.89 | 93.96 |
| ANN | 97.68 | 95.86 | 93.65 | 94.37 |
| KNN | 96.73 | 91.85 | 96.73 | 96.18 |
| RNN | 98.66 | 93.47 | 97.87 | 97.25 |
| CNN | 99.45 | 98.89 | 98.67 | 97.98 |

TABLE V. PERFORMANCE OF CLASSIFICATION WITH OPTIMIZED FEATURES USING NSL-KDD DATASET

| Methods | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| SVM | 92.17 | 93.95 | 89.20 | 91.66 |
| ANN | 95.38 | 94.36 | 91.51 | 90.17 |
| KNN | 91.20 | 90.51 | 95.47 | 94.68 |
| RNN | 97.54 | 95.34 | 94.93 | 90.50 |
| CNN | 98.13 | 97.36 | 97.07 | 95.34 |

## B. Comparative Analysis

This section provides the comparative analysis of proposed HCFS-BOA based CNN model on the evaluation metrics: precision, accuracy, f1-score, and recall as shown in Table VI. The previous methods namely, SafetyMed, HIDM, and BDSDT are employed to assess the HCFS-BOA based CNN performance. SafetyMed [16] achieves 97.63% accuracy,

98.47% precision, 97% recall, and 97.73% f1-score. HIDM [17] achieves 85% accuracy and BDSDT [18] achieves 99.04% accuracy, 87.31% precision, 82.89% recall, and 83.2% f1-score. When compared with the existing methods, the proposed HCFS-BOA based CNN achieves higher accuracy of 99.45%, precision of 98.89%, 98.67% of recall, and 97.98% of f1-score.

TABLE VI. COMPARATIVE ANALYSIS WITH EXISTING METHODS

| Methods | Dataset | Accuracy (%) | Precision (%) | Recall (%) | F1-score (%) |
|---|---|---|---|---|---|
| SafetyMed [16] | CIC-IDS2017 | 97.63 | 98.47 | 97 | 97.73 |
| HIDM [17] | CIC-IDS2017 | 85 | N/A | N/A | N/A |
| BDSDT [18] | CIC-IDS2017 | 99.04 | 87.31 | 82.89 | 83.2 |
| Proposed HCFS-BOA based CNN | CIC-IDS2017 | 99.45 | 98.89 | 98.67 | 97.98 |

## C. Validation of Real-Time Applications

The NSL-KDD dataset is commonly deployed for intrusion detection in IoT to ensure reliability and security for healthcare systems. This research uses the NSL-KDD dataset for training and validation purposes on real-time applications in the cloud. The NSL-KDD dataset is split into training, testing, and validation in the ratio of 70:15:15. IDS is created to detect the different types of attacks by evaluating system logs, network traffic, and behavioral patterns. Malware attacks, DoS attacks, Cross-Site Scripting (XSS), etc., are different attacks. These types of attacks are performed when the patient information is blocked or stolen by attackers. Therefore, the NSL-KDD dataset is employed for model training purposes to reduce the attacks in real-time healthcare applications.

## D. Discussion

The CIC-IDS-2017 dataset is beneficial for intrusion detection because of its comprehensive representation of realistic traffic network scenarios with different types of attacks and normal activities. It provides a labelled and large-scale dataset that assists the evaluation and enhancement of intrusion detection with enhanced robustness and accuracy. The NSL-KDD dataset is beneficial for intrusion detection as it solves limitations in the original KDD Cup dataset by minimizing redundancy and managing a more balanced distribution of classes. It generates the representation of a more realistic modern traffic network that contains normal behavior and different wider attacks that maximize intrusion detection robustness. By using these two datasets, the proposed approach is analyzed by generic type. Moreover, the advantages of the proposed method and the limitations of existing methods are discussed. The existing methods have some limitations such as the SafetyMed method [16] has no defense mechanism against an attack of AML. Combining various detection systems increases the risk of false negatives and false positives in HIDM [17]. BDSDT [18] isn't effective against web and Bot threats since there are fewer instances of these two attack classes. The proposed HCFS-BOA-based CNN model overcomes the existing models' limitations.

To overcome the problem of AML attack, CFS is used to identify highly informative features for minimizing the risk of adversarial manipulations compared to other algorithms. BOA assists in identifying an optimal subset of features that maximizes detection accuracy and reduces the risk of false positives and false negatives. This is done by focusing on informative features in CFS that assist in enhancing the model's ability to discriminate between various attack classes like web and Bot threat. Combining CFS with BOA enables appropriate features that not only have significant correlations with the target variable but also contribute to the optimal performance of intrusion detection in the healthcare system, in contrast to the other methods. The CNN is deployed to identify and classify intrusion accurately and effectively. New attacks such as web and Bot threat attacks are classified effectively by using CNN. The proposed HCFS-BOA-based CNN achieves a superior accuracy of 99.45% when compared with the existing methods namely, SafetyMed, HIDS, and BDSDT.

## V. CONCLUSION

In this research, the HCFS-BOA based on the CNN model is proposed for intrusion detection to secure the entire network in the healthcare system. The proposed method mainly comprises four stages: dataset, min-max normalization, feature selection, and classification. Initially, the data is obtained from the CIC-IDS2017 and NSL-KDD datasets, after which the min-max normalization is performed to normalize the acquired data. For feature selection, HCFS-BOA is employed for optimal performance of intrusion detection in healthcare systems. Finally, the CNN is deployed to identify and classify intrusion accurately and effectively. The proposed HCFS-BOA-based CNN achieves a better accuracy of 99.45% when compared with the existing methods like SafetyMed, HIDS, and BDSDT. In the future, hyperparameter tuning can be applied in feature selection for improving the model's performance.

## REFERENCES

[1] K. Samunnisa, G. S. V. Kumar, and K. Madhavi, "Intrusion detection system in distributed cloud computing: Hybrid clustering and classification methods," Meas.: Sens., vol. 25, p. 100612, Feb. 2023.

[2] A. H. Janabi, T. Kanakis, and M. Johnson, "Overhead Reduction Technique for Software-Defined Network Based Intrusion Detection Systems," IEEE Access, vol. 10, pp. 66481-66491, Jun. 2022.

[3] P. B. Udas, M. E. Karim, and K. S. Roy, "SPIDER: A shallow PCA based network intrusion detection system with enhanced recurrent neural networks," J. King Saud Univ. Comput. Inf. Sci., vol. 34, no. 10B, pp. 10246–10272, Nov. 2022.

[4] O. A. Alzubi, J. A. Alzubi, M. Alazab, A. Alrabea, A. Awajan, and I. Qiqieh, "Optimized Machine Learning-Based Intrusion Detection System for Fog and Edge Computing Environment," Electronics, vol. 11, no. 19, p. 3007, Sep. 2022.

[5] M. Bakro, R. R. Kumar, A. A. Alabrah, Z. Ashraf, S. K. Bisoy, N. Parveen, S. Khawatmi, and A. Abdelsalam, "Efficient Intrusion Detection System in the Cloud Using Fusion Feature Selection Approaches and an Ensemble Classifier," Electronics, vol. 12, no. 11, p. 2427, May 2023.

[6] H. Lin, Q. Xue, J. Feng, and D. Bai, "Internet of things intrusion detection model and algorithm based on cloud computing and multi-feature extraction extreme learning machine," Digital Commun. Networks, vol. 9, no. 1, pp. 111–124, Feb. 2023.

[7] A. K. Samha, N. Malik, D. Sharma, S. Kavitha, and P. Dutta, "Intrusion Detection System Using Hybrid Convolutional Neural Network," Mobile Networks Appl., Aug. 2023.

[8] V. Hnamte, H. Nhung-Nguyen, J. Hussain, and Y. Hwa-Kim, "A Novel Two-Stage Deep Learning Model for Network Intrusion Detection: LSTM-AE," IEEE Access, vol. 11, pp. 37131-37148, Apr. 2023.

[9] R. A. Bakar, X. Huang, M. S. Javed, S. Hussain, and M. F. Majeed, "An Intelligent Agent-Based Detection System for DDoS Attacks Using Automatic Feature Extraction and Selection," Sensors, vol. 23, no. 6, p. 3333, Mar. 2023.

[10] R. Zhao, Y. Mu, L. Zou, and X. Wen, "A Hybrid Intrusion Detection System Based on Feature Selection and Weighted Stacking Classifier," IEEE Access, vol. 10, pp. 71414-71426, Jun. 2022.

[11] R. Shrestha, A. Omidkar, S. A. Roudi, R. Abbas, and S. Kim, "Machine-Learning-Enabled Intrusion Detection System for Cellular Connected UAV Networks," Electronics, vol. 10, no. 13, p. 1549, Jun. 2021.

[12] D. Javeed, M. S. Saeed, I. Ahmad, P. Kumar, A. Jolfaei, and M. Tahir, "An Intelligent Intrusion Detection System for Smart Consumer Electronics Network," IEEE Trans. Consum. Electron., May 2023.

[13] A. Bhardwaj, R. Tyagi, N. Sharma, A. Khare, M. S. Punia, and V. K. Garg, "Network intrusion detection in software defined networking with self-organized constraint-based intelligent learning framework," Meas.: Sens., vol. 24, p. 100580, Dec. 2022.

[14] V. Hnamte and J. Hussain, "Dependable intrusion detection system using deep convolutional neural network: A Novel framework and performance evaluation approach," Telematics and Informatics Reports, vol. 11, p. 100077, Sep. 2023.

[15] S. Shitharth, P. R. Kshirsagar, P. K. Balachandran, K. H. Alyoubi, and A. O. Khadidos, "An Innovative Perceptual Pigeon Galvanized Optimization (PPGO) Based Likelihood Naïve Bayes (LNB) Classification Approach for Network Intrusion Detection System," IEEE Access, vol. 10, pp. 46424-46441, May 2022.

[16] N. Faruqui, M. A. Yousuf, M. Whaiduzzaman, A. K. M. Azad, S. A. Alyami, P. Liò, M. A. Kabir, and M. A. Moni, "SafetyMed: A Novel IoMT Intrusion Detection System Using CNN-LSTM Hybridization," Electronics, vol. 12, no. 17, p. 3541, Aug. 2023.

[17] L. K. Vashishtha, A. P. Singh, and K. Chatterjee, "HIDM: A hybrid intrusion detection model for cloud based systems," Wireless Pers. Commun., vol. 128, no. 4, pp. 2637–2666, Feb. 2023.

[18] P. Kumar, R. Kumar, G. P. Gupta, R. Tripathi, A. Jolfaei, and A. K. M. N. Islam, "A blockchain-orchestrated deep learning approach for secure data transmission in IoT-enabled healthcare system," J. Parallel Distrib. Comput., vol. 172, pp. 69–83, Feb. 2023.

[19] A. Halbouni, T. S. Gunawan, M. H. Habaebi, M. Halbouni, M. Kartiwi, and R. Ahmad, "CNN-LSTM: Hybrid Deep Neural Network for Network Intrusion Detection System," IEEE Access, vol. 10, pp. 99837-99849, Sep. 2022.

[20] H. Han, H. Kim, and Y. Kim, "An Efficient Hyperparameter Control Method for a Network Intrusion Detection System Based on Proximal Policy Optimization," Symmetry, vol. 14, no. 1, p. 161, Jan. 2022.

[21] M. Bakro, R. R. Kumar, A. Alabrah, Z. Ashraf, M. N. Ahmed, M. Shameem, and A. Abdelsalam, "An Improved Design for a Cloud Intrusion Detection System Using Hybrid Features Selection Approach With ML Classifier," IEEE Access, vol. 11, pp. 64228-64247, Jun. 2023.

[22] K. M. Sudar, M. Beulah, P. Deepalakshmi, P. Nagaraj, and P. Chinnasamy, "Detection of Distributed Denial of Service Attacks in SDN using Machine learning techniques," in 2021 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, India, IEEE, 2021, pp. 1-5.

[23] V. Praveena, A. Vijayaraj, P. Chinnasamy, I. Ali, R. Alroobaea, S. Y. Alyahyan, and M. A. Raza, "Optimal deep reinforcement learning for intrusion detection in UAVs," Computers, Materials & Continua, vol. 70, no. 2, pp. 2639-2653, 2022.

[24] P. Chinnasamy, S. Devika, V. Balaji, S. Dhanasekaran, B. J. A. Jebamani, and A. Kiran, "BDDoS: Blocking Distributed Denial of Service Flooding Attacks With Dynamic Path Detectors," in 2023 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, India, IEEE, 2023, pp. 1-5.

[25] P. Chinnasamy, N. Kumaresan, R. Selvaraj, S. Dhanasekaran, K. Ramprathap, and S. Boddu, "An Efficient Phishing Attack Detection using Machine Learning Algorithms," in 2022 International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC), Bhubaneswar, India, IEEE, 2022, pp. 1-6.

[26] E. Anupriya, N. Kumaresan, V. Suresh, S. Dhanasekaran, K. Ramprathap, and P. Chinnasamy, " Fraud Account Detection on Social Network using Machine Learning Techniques," in 2022 International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC), IEEE, 2022, pp. 1-4.

[27] Z. Wu, H. Zhang, P. Wang, and Z. Sun, "RTIDS: A Robust Transformer-Based Approach for Intrusion Detection System," IEEE Access, vol. 10, pp. 64375–64387, Jun. 2022.

[28] S. M. Kasongo, "A deep learning technique for intrusion detection system using a Recurrent Neural Networks based framework," Comput. Commun., vol. 199, pp. 113–125, Feb. 2023.

[29] H. Zhang, B. Zhang, L. Huang, Z. Zhang, and H. Huang, "An Efficient Two-Stage Network Intrusion Detection System in the Internet of Things," Information, vol. 14, no. 2, p. 77, Jan. 2023.

[30] M. A. Siddiqi and W. Pak, "Tier-Based Optimization for Synthesized Network Intrusion Detection System," IEEE Access, vol. 10, pp. 108530–108544, Oct. 2022.

[31] V. Ravuri and S. Vasundra, "Moth-flame optimization-bat optimization: Map-reduce framework for big data clustering using the Moth-flame bat optimization and sparse Fuzzy C-means," Big Data, vol. 8, no. 3, pp. 203-217, Jun. 2020.

[32] S. Vasundra and G. Rajeswarappa, "Hybrid Grasshopper and Improved Bat Optimization Algorithms-based clustering scheme for maximizing lifetime of Wireless Sensor Networks (WSNs)," International Journal of Intelligent Engineering and Systems, vol. 15, pp. 536-546, May. 2022.

[33] S. Vasundra and A. Balaram, "A Hybrid Soft Computing Technique for Software Fault Prediction based on Optimal Feature Extraction and Classification," International Journal of Computer Science and Network Security, vol. 22, no. 5, pp. 348-358, May. 2022.

[34] S. Vasundra and Vasavi Ravuri, "An effective weather forecasting method using a deep long-short-term memory network based on time-series data with sparse fuzzy c-means clustering," Engineering Optimization, vol. 55, no. 9, pp. 1437-1455, Sep. 2022.

# Context-Aware Transfer Learning Approach to Detect Informative Social Media Content for Disaster Management

Saima Saleem*, Monica Mehrotra

Department of Computer Science, Jamia Millia Islamia, New Delhi, India

*Abstract*—In the wake of disasters, timely access to accurate information about on-the-ground situation is crucial for effective disaster response. In this regard, social media (SM) like Twitter have emerged as an invaluable source of real-time user-generated data during such events. However, accurately detecting informative content from large amounts of unstructured user-generated data under such time-sensitive circumstances remains a challenging task. Existing methods predominantly rely on non-contextual language models, which fail to accurately capture the intricate context and linguistic nuances within the disaster-related tweets. While some recent studies have explored context-aware methods, they are based on computationally demanding transformer architectures. To strike a balance between effectiveness and computational efficiency, this study introduces a new context-aware transfer learning approach based on DistilBERT for the accurate detection of disaster related informative content on SM. Our novel approach integrates DistilBERT with a Feed Forward Neural Network (FFNN) and involves multistage finetuning of the model on balanced benchmark real-world disaster datasets. The integration of DistilBERT with an FFNN provides a simple and computationally efficient architecture, while the multistage finetuning facilitates a deeper adaptation of the model to the disaster domain, resulting in improved performance. Our proposed model delivers significant improvements compared to the state-of-the-art (SOTA) methods. This suggests that our model not only addresses the computational challenges but also enhances the contextual understanding, making it a promising advancement for accurate and efficient disaster-related informative content detection on SM platforms.

*Keywords—Disaster management; twitter; distilBERT; deep learning; multistage finetuning; transfer learning*

## I. INTRODUCTION

According to the "Human cost of disasters 2000-2019" report by the "United Nations Office for Disaster Risk Reduction" (UNDRR)[1], disasters have significantly surged over recent decades. As a result, approximately 1.23 million lives have been lost globally, along with an additional economic loss of 2.97 trillion dollars. To effectively respond and manage such significant ecological disruptions, disaster response organizations increasingly rely on swift and precise human-centric information [1]. In recent years, social media (SM) particularly Twitter, have emerged as invaluable tools for

disseminating and obtaining real-time user-generated data during such events [2].

Numerous research works [3-5] have demonstrated that tweets shared on Twitter during disasters often contain informative content, including details about affected people, infrastructure damage, resource needs and availability etc. Such content can be useful for disaster responders in coordinating response efforts, if processed effectively. However, identifying informative content from a substantial volume of irrelevant and noisy tweets during time-critical disaster situations is a challenging task [6]. Moreover, the character limit, uncommon abbreviations, and grammatical mistakes make the detection of informative content even more challenging [7].

Prior research works have explored classical machine learning-based techniques for the automated analysis of SM texts in the context of disasters [8, 9]. However, in recent years, there has been a notable shift towards the application of deep learning (DL) in the analysis of SM text for disaster response [7]. DL has evinced great performance across diverse domains [10-12]. Different DL techniques, including CNN [13, 14], LSTM [15], and Bi-LSTM [16] have been explored for disaster-related tweet classification.

While significant work has been done to analyze SM data related to disasters, the majority of the existing studies rely on traditional non-contextual models within natural language processing (NLP), such as GloVe, Word2Vec (W2V), Bag-of-words (BoW), etc. These models process text sequences in a unidirectional manner, which limits their ability to accurately capture the nuanced context of a tweet sequence. This limitation can lead to inaccurate classification, especially in disaster situations where words like "fire," "flood," and "earthquake" may be used metaphorically. This underscores the need for advanced models capable of understanding the context of disaster-related tweets for the accurate detection of informative content.

In the most recent advancements within NLP, researchers have introduced a spectrum of transfer learning methods and models, notably BERT [17], RoBERTa [18], and DistilBERT [19], among others. These models have demonstrated substantial enhancements across various NLP tasks [20]. They are based on the approach of bidirectional training of transformer-based neural networks to learn the contextual numeric representation of text sequences. This enables them to capture the complete context of a given text sequence, resulting

---

in superior performance compared to conventional non-contextual models [21]. However, they often require significant computational resources due to their large size and complexity.

In the context of disaster response, where speed and efficiency are paramount, the DistilBERT transformer emerges as a promising choice among others. It offers significant advantages in terms of computational efficiency due to its smaller size and faster processing. Despite its demonstrated performance and computational efficiency in various real-time applications like fraud detection [22] and network traffic classification [23], the application of DistilBERT to time-critical disaster management applications remain largely unexplored.

This study aims to fill this research gap by harnessing the capabilities of DistilBERT to offer an effective and computationally efficient solution for the specific task of detecting disaster-related informative content on SM. Our novel approach integrates DistilBERT with a Feed Forward Neural Network (FFNN) and employs multistage finetuning of the model on balanced benchmark datasets of real-world disaster data. The integration of DistilBERT with FFNN provides a simple and computationally efficient architecture, while the multistage finetuning enables a deeper adaptation of the model to the disaster domain, leading to demonstrably improved performance. The contributions of this study are outlined below:

- A new context-aware transfer learning approach is proposed that integrates DistilBERT with a FFNN and involves multistage finetuning of the model to adapt it to the disaster domain for enhanced detection of disaster-related informative content.

- Different model variants are designed using DistilBERT and several DL models like CNN, LSTM, and Bi-LSTM to evaluate the proposed model.

- This study addresses the issue of data imbalance, by employing random down sampling technique to balance the distribution of classes. This ensures that the model is not biased towards the majority class, thus providing unbiased results.

- Through comprehensive ablation studies, this study systematically investigates the impact of multistage finetuning, context-aware representation of tweets, and data balancing on the proposed model's performance.

- The proposed model is compared with various state-of-the-art (SOTA) baseline methods including non-contextual and context-aware transformer-based methods.

- The remaining sections of this paper are arranged as follows: The related works are covered in Section II. Section III covers the methodology. Section IV includes the experimental set up. Section V reports the experimental results. Section VI provides a discussion of the results. Lastly, Section VII concludes this study with the future work.

## II. RELATED WORK

Over the years, extensive research has been carried out and numerous methods have been proposed to extract disaster-related information useful for humanitarian organizations from SM. We group the existing studies into two categories: non-contextual methods and context-aware methods.

### A. Non-contextual Methods

Preliminary studies in this field have mostly employed traditional NLP techniques like BoW along with classical ML algorithms to classify disaster-related SM content. In study [9], the authors presented a system called Tweedr to identify tweets mentioning damage or human fatalities information using a Logistic Regression classifier with BoW features. The authors in [24] extracted situational awareness information from both Hindi and English tweets using lexical and syntactic features with an SVM classifier. In another work [25], the authors employed an SVM classifier along with BoW features to detect tweets related to floods during Manila flooding.

Researchers have also focused on the detection of various information categories present in disaster tweets. In study [26], the authors trained a Naïve Bayes classifier on unigram, bigram, POS, and binary features using the Joplin 2011 tornado disaster dataset to classify tweets into "Caution and advice", "Fatality", "Injury", "Offers of Help", "Missing" and "General Population Information". In another work, a system called AIDR [8] has been developed for identifying informative tweets during disasters. This system extracts features based on BoW model from the tweets and then classifies them into user defined categories such as 'donations,' 'damage' etc.

Additionally, deep neural networks with different word embedding models have been utilized in disaster domain utilizing SM datasets. The authors in study [13, 14] presented CNN-based systems for identifying disaster related SM posts during disasters. They utilized general and domain specific word embedding models for generating features. [15] designed a two-layer LSTM model combined with pretrained GloVe word embeddings for classifying tweet texts of seven different disaster events into informative and not-informative binary categories. In another work [16], the authors employed a Bi-LSTM with pre-trained GloVe embeddings for classifying disaster-related SM textual content.

### B. Context-Aware Methods

With advancements in NLP, researchers have begun to employ pretrained transformer-based architectures in disaster domain to enhance the detection of useful disaster-related information from SM. The authors in study [27] finetuned the BERT model for adapting it to the task of identifying informative tweet texts during disasters. In another work [28], the authors provided a RoBERTa based method, an extension of BERT, for identifying disaster related informative tweet texts. They fine-tuned the architecture to adapt the model to the disaster-specific task.

From the aforementioned studies, it is evident that prevailing methods have majorly relied on traditional non-contextual language models for processing disaster-related SM texts. While a limited number of recent studies have explored

context-aware transformer-based approaches, they typically employ resource-intensive models. Besides imbalanced datasets have been used, which may result in biased outcomes.

## III. METHODOLOGY

Detecting disaster-related informative content on SM effectively and efficiently is crucial for effective disaster response. The problem is framed as a binary classification task: given a tweet, detect whether it is an informative or not_informative tweet. An informative tweet offers valuable details such as information about affected individuals, the extent of infrastructure damage, and resource requirements. Conversely, a non-informative tweet lacks crucial information related to humanitarian organizations or victims. We propose a new context-aware transfer learning approach leveraging the computational efficiency of the DistilBERT language model. Our methodology integrates DistilBERT with a FFNN and employs a multistage finetuning process for improved detection performance. We first provide an overview of the DistilBERT language model and subsequently present our proposed model (DistilBERT+FFNN), followed by multistage finetuning process in detail.

### A. DistilBERT

DistilBERT is a distilled variant of BERT, an advanced bidirectional pretrained language model based on transformer encoded architecture. It is pretrained on a sizable dataset made up of the Wikipedia and Toronto Book Corpus. The reason for choosing DistilBERT is that it is lightweight than the BERT model in terms of parameters and runs 60% faster than BERT. The description of DistilBERT parameters is shown in Table I.

TABLE I.        DISTILBERT PARAMETERS

| Parameter Name | Value |
|---|---|
| Number of Layers | 6 |
| Hidden States Size | 768 |
| Attention Heads | 12 |
| Number of Parameters | 66 million |

DistilBERT comprises six stacked encoders, enabling it to encode the semantic and syntactic information in the tweet sequences as, shown in Fig. 1. The DistilBERT implements a multi-headed self-attention mechanism that captures information from each attention head. This enables the model to look at all the surrounding words in the input tweet sequence, allowing for a better understanding of a word in a particular context.



Fig. 1.   DistilBERT architecture.

This stands in contrast to other text processing models like W2V and GloVe, which generates context-independent embeddings by processing the text sequence in a unidirectional manner. Thus, it can represent each word with a single vector regardless of contextual variations.

### B. Proposed Model

The proposed model is an amalgamation of DistilBERT and an FFNN, leveraging their complementary capabilities to enhance the performance of disaster-related informative content detection task. The DistilBERT is utilized to extract contextual numeric representation of disaster tweet sequences, and the FFNN refines these representations into high-level abstract features for the final prediction. Importantly, FFNN introduces a layer of simplicity to the architecture, ensuring that the model retains its computational efficiency, while maintaining high performance. Fig. 2 illustrates the complete model architecture, showcasing the flow from DistilBERT to FFNN.

To obtain context-aware vector representation of the preprocessed tweets from DistilBERT model, tweets must be transformed into a format understandable by DistilBERT. For this, the tweets are passed to the batch_encode_plus method of DistilBertTokenizer class from the transformers package. It performs the following operations to transform the textual data into an appropriate format:

- Tokenization: Splitting the tweet text sequence into a list of words or sub-words called tokens.

- Padding: Making the length of the tweet sequences equal in case of the unequal sequence lengths.

- Adding special tokens: Adding special tokens such as [CLS], which stand for classification, and [SEP], which stands for separation, to indicate the beginning and the end of each sentence, respectively.

- Encoding: Substituting tokens with their corresponding IDS.

- Adding attention mask: Including an attention mask, a binary array guiding the model on which tokens to focus on it and which to ignore.

For each tweet, the *batch_encode_plus* method returns two sequences (input IDs along with attention mask), which are then input to the DistilBERT. The DistilBERT model outputs hidden states of shape (batch_size, sequence_length, hidden_size), representing the word-level/token-level embedding output of DistilBERT's layers. In this study, the output of the last hidden state is considered as it typically leads to the best empirical results [29]. Moreover, instead of using the word-level/token-level representation, the sentence-level representation of the sequence is used by taking the output for the [CLS] token, denoted as $R_{[CLS]}$. The sentence-level embedding $R_{[CLS]}$ provides the overall context of the entire sequence of tweet text. $R_{[CLS]}$ is a 2D tensor of shape (batch_size, hidden_size), which is passed to the next component of our model i.e., FFNN for predicting whether a tweet is informative or not-informative.

Fig. 2. Proposed model (DistilBERT+FFNN) for the detection of disaster-related informative tweets.

The FFNN is comprised of two dense layers: a hidden layer and an output layer. The hidden layer uses a popular Rectified Linear Unit (ReLU) activation function. The ReLU function is computed as shown in Eq. (1) and Eq. (2):

$$(x) = \max(0, x) \tag{1}$$

$$f(x) = \begin{cases} 0, & x < 0 \\ x, & x > 0 \end{cases} \tag{2}$$

It outputs 0 for every value of $x < 0$ and $x$ itself for all values of $x>0$. The output of the FFNN is computed as shown in the formula in Eq. (3) and Eq. (4):

$$Z = \sum_{i=0}^{n} W_i \times X_i + B \tag{3}$$

$$f(z) = f\left(\sum_{i=0}^{n} W_i \times X_i + B\right) \tag{4}$$

where, $f$ is an activation function (i.e., sigmoid for the output layer), X is the previous layer output, W is the weight matrix, n is the number of inputs from the incoming layer, and B is the bias.

The output of the sigmoid function is always between 0 and 1. The sigmoid activation function is computed in Eq. (5).

$$y' = \frac{1}{1+e^{-z}} \tag{5}$$

where, y' is the model's predicted value, z is the output generated by the last layer of FFNN as computed in (3). Since, the number of classes is two in our case, Binary Cross Entropy (BCE) loss function is employed which computes how far the model deviates from the correct prediction i.e., error as illustrated in Eq. (6).

$$L_{BCE} = -\frac{1}{N} \sum_{i=1}^{N} (y_i \cdot log y'_i + (1-y_i) \cdot \log(1-y'_i)) \tag{6}$$

where N is the training set size, $y_i$ and $y_i'$ is the actual class label and predicted value for the $i^{th}$ sample in the dataset. BCE outputs a loss value that tells how wrong the models' predictions are. The lower the loss value, the higher the model's performance in making accurate predictions.

### C. Multistage Finetuning Procedure

This study employs a multistage finetuning approach to finetune the proposed model. In this approach, the proposed model undergoes finetuning in a series of two stages. In the first stage (stage-1), the DistilBERT model is frozen and only the FFNN undergoes finetuning. This allows the FFNN to learn the task-specific knowledge without altering the general knowledge acquired by DistilBERT during pretraining. The model is trained for six epochs using Adam optimizer with a learning rate of 5e-5 and batch-size of 64. During the error

back-propagation, the pretrained weights of DistilBERT will remain unchanged; only the FFNN will learn. Once, the weights of FFNN are learned in the subsequent stages (stage-2), the DistilBERT layers are unfrozen, and the entire architecture undergoes further finetuning for six additional epochs. The pretrained DistilBERT weights also get updated during finetuning, implying that the error gets back-propagated through the entire architecture. A lower learning rate is set to prevent significant updates to the gradient. A callback mechanism is used to stop the training process when the model has stopped improving. The DistilBERT attention dropout and DistilBERT dropout are also slightly increased from their default values. This two stage finetuning process allows the model to further adapt to our task and improve overall performance.

## IV. EXPERIMENTAL SETUP

This section presents the setup for experiments, including datasets used, evaluation metrics, model variants, baseline methods, and training details.

### A. Dataset Description and Pre-Processing

The current study evaluates the proposed model on a real-world disaster tweet dataset called CrisisMMD [30]. This dataset encompasses tweets from seven disaster events broadly annotated into informative and not_informative classes. A few examples of informative and not_informative tweets from the CrisisMMD dataset are shown in Fig. 3. An overview of the number of tweets present in each disaster dataset is provided in Table II and the class distribution of each disaster dataset can be seen in Fig. 4.

Before conducting any experiments, the datasets are cleaned by eliminating duplicate tweets, hashtags, URLs, user mentions and various symbols such as "@," "!","#", "&," and "%". As can be observed from Fig. 4, the dataset is imbalanced, to address the class imbalance and ensure equal representation of *informative* and *not_informative* tweets, a balanced sample is obtained using random down-sampling. This simple and effective technique reduces the majority class to match the size of the minority class. Finally, a total count of about 8.6k tweet samples from the CrisisMMD dataset are used for subsequent processing. From each event dataset 70% tweet samples are used for training, 10% are used for validation, and the remaining 20% are used for testing the model's performance.

TABLE II. CRISISMMD DATASET DETAILS

| Disaster event | #Tweets |
|---|---|
| Hurricane Harvey | 4434 |
| Hurricane Maria | 4556 |
| Hurricane Irma | 4504 |
| Sri Lanka Floods | 1022 |
| Iraq-Iran Earthquake | 597 |
| Mexico Earthquake | 1380 |
| California Wildfires | 1588 |



Fig. 3. Examples of informative and not_informative tweets from the CrisisMMD dataset.



Fig. 4. CrisisMMD class distribution.

### B. Evaluation Metrics

To assess our proposed model's performance, several metrics are used:

*1) Precision (P):* shows the proportion of correctly predicted informative tweets to the total predicted informative tweets and is equal to:

$$P = \frac{TP}{TP+FP} \tag{7}$$

*2) Recall (R):* shows the proportion of correctly predicted informative tweets to the total actual informative tweets and is equal to:

$$R = \frac{TP}{TP+FN} \tag{8}$$

*3) F1-Score (F1):* merges the precision and recall by computing their harmonic mean, as:

$$F1 = 2 \times \frac{P \times R}{P+R} \tag{9}$$

where, TP indicates "True Positive", which means an informative tweet is correctly predicted as informative. TN indicates "True Negative", which means a not_informative tweet is correctly predicted as not_informative. FP indicates "False Positive", which means a not_informative tweet is incorrectly predicted as informative, and FN indicates "False Negative", which means an informative tweet is incorrectly predicted as not_informative.

*C. Model Variants*

To evaluate the effectiveness of the proposed model (DistilBERT+FFNN), we design three different model variants. These variants are formed by integrating DistilBERT with various popular DL architectures specifically CNN, LSTM, and Bi-LSTM. The description of each model variant is provided below.

*1) DistilBERT+CNN:* CNN is combined on top of DistilBERT as a classification component and the input to CNN is a 2-dimensional tensor *X* obtained from DistilBERT. This tensor undergoes a convolution operation with a filter matrix, generating a new feature map through element-wise multiplication. The resulting feature map is then subjected to a pooling layer, extracting maximum values to form a pooled output. This output is subsequently fed into the output layer, to determine the tweet's class as shown in Eq. (4).

*2) DistilBERT+LSTM:* An LSTM is used on top of DistilBERT and the input to the LSTM is DistilBERT output. An LSTM comprises recurrently connected memory units. Each unit comprises of a cell state, an input gate, an output gate and a forget gate. The cell state keeps information over arbitrary periods, and the information flow into and out of the cell state is governed by the gates. This allows the model to retain important information. The output from the LSTM layer is fed to an output layer for classification as shown in Eq. (4).

*3) DistilBERT+Bi-LSTM:* Another variant is designed by using a Bi-LSTM on top of DistilBERT. A Bi-LSTM layer trains two separate LSTM layers of opposite directions (forward and backward) simultaneously on the input generated from DistilBERT and then concatenates the outputs from both layers. This output is fed to a final output dense layer for classifying tweets into informative and not-informative classes as shown in Eq. (4).

*D. Baseline Methods*

The proposed model is further evaluated by comparing it against the SOTA baseline methods for disaster-related informative content detection. These methods are grouped into two categories: non-contextual methods and context-aware methods. A brief description of each baseline method is provided below:

*1) Non-contextual methods:*

*a)* W2V-CNN [13]: The authors in [13] employed a CNN model with filters of different sizes for identifying disaster related SM posts. They used general pretrained W2V embeddings with the CNN model.

*b)* CW2V-CNN [14]: The authors in [14] used a pretrained crisis embedding model (CW2V) [31] and trained a custom CNN with different filters to identify disaster related SM posts during a disaster.

*c)* GloVe-LSTM [15]: The authors in [15] used a pretrained GloVe word embeddings of 100 dimension with a 2-layer LSTM model for identifying informative textual content from SM during disasters.

*d)* GloVe-Bi-LSTM [16]: The authors in [16] used a pretrained GloVe word embeddings of 300 dimension along with a Bi-LSTM neural network for detecting informative textual content from SM during disasters.

*2) Context-aware methods:*

*a)* Finetuned-BERT [27]: The authors in [27] finetuned BERT for classifying SM textual posts into *informative* and *not_informative* classes during disasters.

*b)* Finetuned-RoBERTa [28]: The authors in [28] used a variant of BERT called RoBERTa and finetuned it for identifying informative SM textual posts during disasters.

All of these baseline methods have used the CrisisMMD dataset, with the exception of W2V-CNN [13], CW2V-CNN [14]. For these specific methods, we rely on results computed in [32] on the same CrisisMMD dataset, to facilitate a comprehensive and consistent comparison in this study.

*E. Training Details*

All the experiments are executed using TensorFlow framework and Google Colab cloud platform with Python programming language. The optimal hyperparameters for the stage-1 and stage-2 of our model finetuning approach are listed in Table III and IV, respectively.

TABLE III.    PARAMETERS FOR STAGE-1 OF FINETUNING

| Hyperparameter | Value |
|---|---|
| Learning rate | 5e-5 |
| Number of epochs | 6 |
| Batch Size | 64 |
| Dropout | 0.2 |

TABLE IV.    PARAMETERS FOR STAGE-2 OF FINETUNING

| Hyperparameter | Value |
|---|---|
| Learning rate | 2e-5 |
| Number of epochs | 6 |
| Batch Size | 64 |
| DistilBERT Dropout | 0.2 |
| DistilBERT Attention Dropout | 0.2 |

V.    EXPERIMENTAL RESULTS

In this section, we present a comprehensive evaluation of the performance of the proposed model and its various variants employed in this study. Additionally, a thorough comparison is conducted between the proposed model and SOTA baseline methods. Finally, this section covers detailed ablation studies.

## A. Proposed Model vs. Different Model Variants

To demonstrate that the proposed model (DistilBERT+FFNN) is an effective model for the informativeness classification task, it is compared against different model variants, as discussed in Section C. The P, R, and F1 of the proposed model and the model variants on CrisisMMD dataset are reported in Table V, with bold values indicating the best results. As per the table, all the model variants exhibit good performance across seven disasters. The proposed model exhibits the highest P ranging from 73.12 to 96.20, R in the range of 73.11 to 96.04, and F1 in the range of 72.20 to 96.04. These results render the proposed model an effective choice for detecting disaster related informative content.

## B. Proposed Model vs. SOTA Baseline Methods

In this subsection, we first present a performance comparison of the proposed model against non-contextual baseline methods. Next, we analyze the comparative performance of the proposed model with respect to context-aware baseline methods.

*1) Proposed model vs. non-contextual baseline methods:* The comparison of results in terms of F1 of the proposed model against non-contextual methods: W2V-CNN [13], CW2V-CNN [14], GloVe-LSTM [15] and GloVe-Bi-LSTM [16] are reported in Table VI. As per the results, the proposed model consistently outperforms the non-contextual methods across all disasters, with the exception of Hurricane Irma, where GloVe-LSTM [15] outperforms in terms of F1. On average, the proposed model significantly enhances F1, with an improvement ranging from 2.75% to 19.86%.

*2) Proposed model vs context-aware baseline methods:* The efficacy of the proposed model is further validated through comparisons with recent transformer-based context-aware methods: Finetuned-BERT [27] and Finetuned-RoBERTa [28]. The corresponding results, expressed in terms of F1, are detailed in Table VII. A comprehensive examination of the table reveals that the proposed model outperforms Finetuned-BERT [27] across six out of seven disasters and surpasses Finetuned-RoBERTa [28] on five out of seven disasters from the CrisisMMD dataset. The proposed model exhibits an average F1 improvement of 5.47% over Finetuned-BERT [27] and 1.6% over Finetuned-RoBERTa [28].

## C. Ablation Studies

In this subsection, we conduct ablation experiments to investigate the effect of multistage finetuning, context-aware representation of tweets, and data balancing on the proposed model performance.

*1) Effect of multistage finetuning:* To understand the effect of multistage finetuning on the model performance, we conduct a comparative analysis with a single stage (stage-1) finetuned model, where only the FFNN is finetuned while keeping DistilBERT parameters frozen. The results depicted in Fig. 5 report the F1 achieved with and without multistage finetuning on CrisisMMD dataset. From the results, consistent improvement is observed with multistage finetuning, where the model undergoes a series of two finetuning stages. Specifically, multistage finetuning boosts the performance by 3.89% to 6.85% in terms of F1. This emphasizes the efficacy of multistage finetuning in enhancing the model's ability to capture and adapt to the nuances present in disaster-related tweets, resulting in improved performance.

*2) Effect of DistilBERT context-aware tweet representation:* To evaluate the impact of context-aware tweet representation derived from DistilBERT on the informativeness classification task, additional experiments are performed, comparing DistilBERT against two non-contextual word embedding models. Experiments are designed where the FFNN classifier uses CW2V [33] and GloVe [34] embedding models.

TABLE V.        COMPARISON OF THE PROPOSED MODEL VS. DIFFERENT VARIANTS

| Disasters | Models | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DistilBERT+CNN | | | DistilBERT+LSTM | | | DistilBERT+Bi-LSTM | | | Proposed Model | | |
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| Hurricane Harvey | 81.23 | 81.23 | 81.23 | 81.10 | 81.10 | 81.10 | 83.24 | 83.21 | 83.00 | **84.57** | **84.56** | **84.56** |
| Hurricane Maria | 83.03 | 83.03 | 83.02 | 82.21 | 82.20 | 82.20 | 81.31 | 81.30 | 81.25 | **84.03** | **84.01** | **83.11** |
| Hurricane Irma | 82.00 | 81.22 | 81.19 | 80.05 | 80.05 | 80.04 | 82.04 | 82.04 | 82.01 | **83.12** | **83.07** | **82.12** |
| Sri Lanka Floods | 94.09 | 94.09 | 94.08 | 93.01 | 93.00 | 93.00 | 94.09 | 94.04 | 94.00 | **96.20** | **96.04** | **96.04** |
| Mexico Earthquake | 71.10 | 71.09 | 71.09 | 73.06 | 72.31 | 71.40 | **73.19** | 73.01 | 72.24 | **73.19** | **73.17** | **73.17** |
| Iraq-Iran Earthquake | 81.04 | 81.03 | 81.03 | 80.12 | 80.10 | 80.10 | 81.11 | 81.10 | 81.09 | **82.09** | **82.05** | **82.02** |
| California Wildfires | 72.09 | 71.30 | 71.30 | 72.03 | 72.01 | 72.01 | 72.06 | 72.06 | 72.03 | **73.12** | **73.11** | **72.20** |

TABLE VI.    COMPARISON OF THE PROPOSED MODEL VS. NON-CONTEXTUAL BASELINE METHODS IN TERMS OF F1

| Disasters | Methods | | | | |
|---|---|---|---|---|---|
| | *W2V-CNN [13]* | *CW2V-CNN [14]* | *GloVE-LSTM [15]* | *GloVe -Bi-LSTM [16]* | *Proposed Model* |
| Hurricane Harvey | 74.20 | 75.00 | 81.00 | 70.93 | **84.56** |
| Hurricane Maria | 69.10 | 70.00 | 82.00 | 69.51 | **83.11** |
| Hurricane Irma | 69.60 | 70.00 | **83.00** | 57.17 | 82.12 |
| Sri Lanka Floods | 90.02 | 91.00 | 92.00 | 72.57 | **96.04** |
| Mexico Earthquake | 64.20 | 64.00 | 71.00 | 51.83 | **73.17** |
| Iraq-Iran Earthquake | 59.10 | 60.00 | 81.00 | 63.43 | **82.02** |
| California Wildfires | 57.40 | 58.00 | 64.00 | 48.81 | **72.20** |
| **Average F1** | 69.09 | 69.71 | 79.14 | 62.03 | **81.89** |

TABLE VII.    COMPARISON OF THE PROPOSED MODEL VS. CONTEXT-AWARE BASELINE METHODS IN TERMS OF F1

| Disasters | Methods | | |
|---|---|---|---|
| | *Finetuned-BERT [27]* | *Finetuned-RoBERTa [28]* | *Proposed Model* |
| Hurricane Harvey | 83.00 | 84.50 | **84.56** |
| Hurricane Maria | 79.00 | **88.00** | 83.11 |
| Hurricane Irma | 73.00 | 82.00 | **82.12** |
| Sri Lanka Floods | 96.00 | 95.50 | **96.04** |
| Mexico Earthquake | **78.00** | 74.00 | 73.17 |
| Iraq-Iran Earthquake | 63.00 | 72.50 | **82.02** |
| California Wildfires | 63.00 | 65.50 | **72.20** |
| **Average F1** | 76.42 | 80.29 | **81.89** |



Fig. 5.    F1 of the proposed model with and without multistage finetuning.

The CW2V is a 300-dimensional embedding, pretrained using W2V model on disaster related Twitter corpus. GloVe is a pretrained 100-dimensional embedding trained on Wikipedia and web text words. Fig. 6 depicts the F1 of the proposed model and non-contextual models. Our investigation reveals that the context-aware representation generated by DistilBERT model improves the results across all disasters compared to the non-contextual models. Overall, with DistilBERT, a significant improvement in the range of 8.11% to 19.56% compared to GloVe and 9.75% to 21.2% compared to CW2V in F1 is achieved across seven disasters. These results clearly highlight that DistilBERT model captures better contextual information from tweet sequences leading to better detection performance than GloVe and CW2V models.



Fig. 6.    F1 of the proposed model and non-contextual models.

*3) Effect of data balancing*: As mentioned earlier, the datasets are imbalanced, and random down sampling is employed to ensure a balanced representation. Table VIII presents a comparative analysis of the proposed model's

performance in terms of F1 when trained on imbalanced and balanced dataset. From the table, it can be concluded that the model's performance is notably enhanced, showcasing improved and unbiased results when trained on balanced datasets. On the other hand, training on imbalanced datasets leads the model to disproportionately learn from the majority class, thereby introducing bias towards that class.

TABLE VIII.   F1 ANALYSIS OF THE PROPOSED MODEL ON IMBALANCED VS. BALANCED DATASET

| Disasters | Class 0 (*not_informative*) 1 (*informative*) | Imbalanced | Balanced |
|---|---|---|---|
| Hurricane Harvey | 0 | 61.51 | **84.69** |
| | 1 | 86.24 | **84.43** |
| Hurricane Maria | 0 | 73.89 | **82.92** |
| | 1 | 89.33 | **83.30** |
| Hurricane Irma | 0 | 61.45 | **82.03** |
| | 1 | 80.24 | **82.21** |
| Sri Lanka Floods | 0 | 94.42 | **95.59** |
| | 1 | 71.09 | **96.49** |
| Mexico Earthquake | 0 | 48.80 | **73.10** |
| | 1 | 79.02 | **73.24** |
| Iraq-Iran Earthquake | 0 | 42.87 | **82.00** |
| | 1 | 79.86 | **82.04** |
| California Wildfires | 0 | 36.58 | **72.10** |
| | 1 | 87.45 | **72.30** |

## VI.   DISCUSSION

Upon a thorough evaluation of our proposed model on seven real-world disasters from a benchmark CrisisMMD dataset, it stands out as both effective and computationally efficient for the informativeness classification task. In the performance comparison between the proposed model and its variants, it is noteworthy that the integration of CNN, LSTM, and Bi-LSTM DL models with DistilBERT does not enhance the classification performance. Conversely, a simple FFNN proves sufficient for extracting high-level abstract features from the contextualized representations generated from DistilBERT. The results confirm the effectiveness of the proposed model for the informativeness classification task. The outperformance of the proposed model against non-contextual methods based on W2V, CW2V and GloVe models underscores the importance of context-aware representation of tweets in the effective detection of disaster-related informative content. The comparison of the proposed model with SOTA context-aware methods based on BERT and RoBERTa, demonstrates clear advantages of the proposed model in the detection performance. While BERT and RoBERTa are powerful transformer architectures renowned for their contextual understanding, require significant computational resources due to their large size and complexity. The proposed model built on DistilBERT known for its enhanced computational efficiency capitalizes on the efficiency gains through the multi-stage fine-tuning process. The ablation experimental results provide deeper insights, emphasizing the

advantages of our multistage finetuning approach for significantly enhancing overall model performance. Additionally, the combination of DistilBERT with an FFNN, and data balancing collectively contributes to the effective and computationally efficient detection of disaster-related informative content on SM. This positions the model as well-suited for time-critical disaster response applications. Nevertheless, one of the limitations of this study is that it considers only English language tweets, while people often communicate in their native language during disasters. So, multilingual transformer-based models are worth investigating to tackle multi-linguality issues.

## VII.   CONCLUSION AND FUTURE WORK

In the realm of disaster management, this study introduces a novel context-aware transfer learning approach harnessing the computational efficiency of DistilBERT for the detection of disaster-related informative content on SM. Our methodology integrates DistilBERT with an FFNN, providing a simple yet effective architecture. A key feature of our approach involves multistage finetuning of the model on seven real-world disasters, resulting in improved detection performance.

This work contributes to the broader goal of improving the effectiveness and efficiency of disaster response systems using NLP and AI technologies. By developing a specialized model for disaster-related informative content detection, we aim to provide a valuable tool for disaster response organizations to better identify critical information during disasters and provide situational awareness to decision-makers. The model can be implemented in any system for filtering actionable informative content during disasters from irrelevant content, enabling disaster responders to improve their ability to deliver help and make well-informed decisions. Furthermore, the proposed model holds promise for diverse domains such as epidemics and civil unrest monitoring on SM. With domain-specific finetuning, it can be readily adapted to identify informative content during outbreaks, protests, or other volatile situations enabling real-time interventions.

In the future, we plan to extend this work for identifying and categorizing multimodal informative content into distinct humanitarian information categories, including "affected individuals", "infrastructure damage", "help and rescue", "resource needs" etc. to enable a more targeted and efficient disaster response. Moreover, we recognize the need to explore cross-domain informative content detection, particularly in situations where labeled data for the ongoing disaster might be scarce or unavailable.

## REFERENCES

[1] Wang, Z., & Ye, X. (2018). Social media analytics for natural disaster management. *International Journal of Geographical Information Science* 32(1), 49-72.

[2] Saleem, S., & Mehrotra, M. (2022). Emergent Use of Artificial Intelligence and Social Media for Disaster Management. In *Proceedings of International Conference on Data Science and Applications: ICDSA 2021*, *Volume 2* (pp. 195-210). Springer Singapore.

[3] Saleem, S., & Mehrotra, M. (2023, July). An Analytical Framework for Analyzing Tweets for Disaster Management: Case Study of Turkey Earthquake 2023. In *2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT)* (pp. 1-7). IEEE.

[4]     Vieweg, S. (2010) Microblogged contributions to the emergency arena: Discovery, interpretation and implications. In *Proceedings of Computer Supported Collaborative Work* (pp. 515-516).

[5]     Alam, F., Alam, T., Hasan, M. A., Hasnat, A., Imran, M., & Ofli, F. (2023). MEDIC: a multi-task learning dataset for disaster image classification. *Neural Computing and Applications*, 5(3), 2609-32.

[6]     Caragea, C., Silvescu, A., & Tapia, A. H. (2016, May). Identifying informative messages in disaster events using convolutional neural networks. In *International conference on information systems for crisis response and management* (pp. 137-147).

[7]     Imran, M., Ofli, F., Caragea, D., & Torralba, A. (2020). Using AI and social media multimodal content for disaster response and management: Opportunities, challenges, and future directions. *Information Processing & Management*, 57(5), 102261.

[8]     Imran, M., Castillo, C., Lucas, J., Meier, P., & Vieweg, S. (2014). AIDR: Artificial intelligence for disaster response. In *Proceedings of the 23rd International Conference on World Wide Web* (pp. 159-162).

[9]     Ashktorab, Z., Brown, C., Nandi, M., & Culotta, A. (2014). Tweedr: Mining twitter to inform disaster response. In *Proceedings of the 11th International ISCRAM Conference* (pp. 269-272).

[10]   Tunio, M. H., Jianping, L., Butt, M. H. F., Memon, I., & Magsi, Y. (2022). Fruit Detection and Segmentation Using Customized Deep Learning Techniques. In *2022 19th International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP)* (pp. 1-5). IEEE.

[11]   Dhaka, D., Saleem, S., & Mehrotra, M. (2023). USE-Based Deep Learning Approach to Detect Spammers on Twitter. In *International Conference on Smart Trends in Computing and Communications* (pp. 407-416). Singapore: Springer Nature Singapore.

[12]   Ajmal, S., Sarfraz, M. S., Memon, I., Bilal, M., & Alam, K. A. (2023). PUB-VEN: a personalized recommendation system for suggesting publication venues. Multimedia Tools and Applications, 1-22.

[13]   Burel, G., & Alani, H. (2018). Crisis event extraction service (crees)- automatic detection and classification of crisis-related content on social media. In *Proceedings of the 15th International Conference on Information Systems for Crisis Response and Management (ISCRAM).*

[14]   Alam, F., Ofli, F., & Imran, M. (2019). CrisisDPS: Crisis Data Processing Services. In *Proceedings of the 16th International Conference on Information Systems for Crisis Response and Management (ISCRAM).*

[15]   Kumar, A., Singh, J., Dwivedi, Y., & Rana, N. (2020). A deep multi-modal neural network for informative Twitter content classification during emergencies. *Annals of Operations Research*, 1-32.

[16]   Khattar, A., Quadri, S. M. K (2022). CAMM: Cross-Attention Multimodal Classification of Disaster-Related Tweets. *IEEE Access*, 10, 92889-902.

[17]   Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805.*

[18]   Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., … & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692.*

[19]   Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108.*

[20]   Kalyan, K. S., Rajasekharan, A., & Sangeetha, S. (2021). Ammus: A survey of transformer-based pretrained models in natural language processing. *arXiv preprint arXiv:2108.05542.*

[21]   Kaliyar, R. K., Goswami, A., & Narang, P. (2021). FakeBERT: Fake news detection in social media with a BERT-based deep learning approach. *Multimedia tools and applications,* 80(8), 11765-88.

[22]   Chang, J. W., Yen, N., & Hung, J. C. (2022). Design of a NLP-empowered finance fraud awareness model: the anti-fraud chatbot for fraud detection and fraud classification as an instance. *Journal of Ambient Intelligence and Humanized Computing,* 4663-79.

[23]   Shin, C. Y., Park, J. T., Baek, U. J., & Kim, M. S. (2023). A Feasible and Explainable Network Traffic Classifier Utilizing DistilBERT. *IEEE Access.*

[24]   Rudra, K., Ganguly, N., Goyal, P., & Ghosh, S. (2018). Extracting and summarizing situational information from the twitter social media during disasters. *ACM Transactions on the Web (TWEB)*, 12(3), 1-35.

[25]   Parilla-Ferrer, B. E., Fernandez, P. L., Ballena, J. T. (2014). Automatic classification of disaster-related tweets. In *Proceedings of the International Conference on innovative engineering technologies (ICIET)* (pp 62-69).

[26]   Imran, M., Elbassuoni, S., Castillo, C., Diaz, F., & Meier, P. (2013). Extracting information nuggets from disaster-Related messages in social media. In *International conference on information systems for crisis response and management* (pp. 791-801).

[27]   Madichetty, S., Muthukumarasamy, S., & Jayadev, P. (2021). Multi-modal classification of Twitter data during disasters for humanitarian response. *Journal of ambient intelligence and humanized computing,* 10223-37.

[28]   Madichetty, S., & Madisetty, S. (2023). A RoBERTa based model for identifying the multi-modal informative tweets during disaster. *Multimedia Tools and Applications,* 1-19.

[29]   Sun, C., Qiu, X., Xu, Y., & Huang, X. (2019). How to fine-tune bert for text classification?. In *Chinese Computational Linguistics: 18th China National Conference, CCL 2019, Kunming, China, October 18–20, 2019, Proceedings 18* (pp. 194-206). Springer International Publishing.

[30]   Alam, F., Ofli, F., & Imran, M. (2018, June). Crisismmd: Multimodal twitter datasets from natural disasters. In *Proceedings of the Twelfth international AAAI conference on web and social media* (Vol. 12, No. 1).

[31]   Alam, F., Joty, S., & Imran, M. (2018, June). Graph based semi-supervised learning with convolution neural networks to classify crisis related tweets. In *Proceedings of the international AAAI conference on web and social media* (Vol. 12, No. 1).

[32]   Madichetty, S. (2020). Classifying informative and non-informative tweets from the twitter by adapting image features during disaster. *Multimedia Tools and Applications*, 79, 28901-28923.

[33]   Nguyen, D. T., Joty, S., Imran, M., Sajjad, H., & Mitra, P. (2016). Applications of online deep learning for crisis response using social media information. *arXiv preprint arXiv:1610.01030.*

[34]   Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543).

# Practical Application of AI and Large Language Models in Software Engineering Education

Vasil Kozov[1], Galina Ivanova[2], Desislava Atanasova[3]

Dept. of Informatics, University of Ruse "Angel Kanchev", Ruse, Bulgaria[1, 3]

Dept of Computer Systems and Technology, University of Ruse "Angel Kanchev", Ruse, Bulgaria[2]

*Abstract*—Subjects with limited application in the software industry like AI have recently received tremendous boon due to the development and raise of publicity of LLMs. LLM-powered software has a wide array of practical applications that must be taught to Software Engineering students, so that they can be relevant in the field. The speed of technological change is extremely fast, and university curriculums must include those changes. Renewing and creating new methodologies and workshops is a difficult task to complete successfully in such a dynamic environment full of cutting-edge technologies. This paper aims to showcase our approach to using LLM-powered software for AI generated images, like Stable diffusion and code generation tools like ChatGPT in workshops for two relevant subjects – Analysis of Software Requirements and Specifications, as well as Artificial Intelligence. A comparison between the different available LLMs that generate images is made, and the choice between them is explained. Student feedback is shown and a general positive and motivational impact is noted during and after the workshop. A brief introduction that covers the subjects where AI is applied is made. The proposed solutions for several uses of AI in the field of higher education, more specifically software engineering, are presented. Several workshops have been made and included in the curriculum. The results of their application have been noted and an analysis is made. More propositions on further development based on the gained experience, feedback and retrieved data are made. Conclusions are made on the application of AI in higher education and different ways to utilize such tools are presented.

*Keywords*—*Application of AI-powered software; AI generated images; software engineering; stable diffusion; higher education*

## I. INTRODUCTION

With the release of popular and easy to use text based large language models, Artificial intelligence (AI) chatbots and image generation AI, it has become a necessity to teach students how to properly use them for the correct purposes. A big detriment of using AI by first and second year students is their lack of understanding in the subjects where they are trying to use AI. This sadly leads to less problem-solving thinking and doing actual effort to work on tasks. This in turn leads to less brain development and lack of practical skills. Teaching students how to correctly use such tools to enhance their learning is imperative for developing them into good specialists that may in turn use AI powered software and hardware to their benefit.

The Software Engineering bachelor's degree was chosen for the inclusion of teaching AI. Based on our experience, the more technically advanced and tech savvy students are more receptive towards the inclusion of new and experimental approaches. This was the main factor in our choice of who to use this approach with. Several subjects were selected, where the methodologies were applied. Workshops and materials to help with training the students were created. In the current report the focus will be on the application of Large Language Models (LLMs) in two subjects - Artificial intelligence (6th semester) and Analyzing system requirements and specifications (5th semester). It will also be discussed how to use LLMs to enhance the subject Introduction to Programming by writing unit tests (while that is still in the testing phase and has only been partially applied), as well as future ideas for the inclusion of image generation AI in Computer Graphics.

In order to apply the tools that contain AI, methodologies have been adapted to suit the needs of the subjects. Two of them will be discussed and the choices that have been made as well as the reasoning behind them will be explained.

## II. APPLICATION OF AI-POWERED SOFTWARE

### A. Inclusion of AI Assistance in Analyzing Software Requirements and Specifications (ASRS)

At this point in time during their university education, software engineering students have already passed Databases (DBs), Object Oriented Programming (OOP), and some of their Algorithm subjects. They are introduced to the vision on how to lighten their workload using assistive tools powered by AI. ChatGPT is the text generation model that is currently being used, but as new AI chat bots emerge, university staff is testing them and showing students the differences between them.

In the subject ASRS, the students are required to submit a project – a software system that they must create and fully document using an iterative process. Our first use of ChatGPT is when we demonstrate its capabilities for idea generation.

As seen on Fig. 1, a description of the problem help is needed with is given on the left, and although it lacks detail, it is a true and valid statement, so the model gives a useful, albeit generic answer. On the right side the LLM is queried using a message that has more details included. It has been decided to include as much relevant information as possible as the specific requirements are narrowed down. A description of the field of study and level of progress in the field is defined – "student in Software engineering", this can be further expanded by specifying that this is a bachelor's degree. By noting the current user skillset, the LLM can get software project ideas within the correct scope, omitting unfamiliar technologies. The

"course project" part of the message is also a component that will be used to further refine the scope of the ideas, giving them a timeframe. After defining what the LLM is required to do, the specified field is narrowed even further by statements of likes and dislikes. The overall structure of the message on the right helps the LLM print ideas that are more likely to interest the user and are within their means to achieve. The term "prompt engineering" [1] can better help describe this part of the communication with AI. It is highly possible that every person working with LLMs will be required to understand how to express themselves in a meaningful way if they want to be more efficient in the long term. Whether or not that will become a job requirement is speculative, but it is certainly a useful skill to have right now for students.

The logical path in narrowing down the choices that the model must make in its answer looks similar to this hierarchy in this case: field (of study or work) -> constraints (skillset, abilities, requirements) -> environment descriptions (question framework) -> task that has to be completed by the LLM (actual answer format) -> formatting (table, software code, bullet points, queries for other software) -> further requirements (inclusion and exclusion of topics and fields).

It is of course allowed for the students to use their own ideas irrespective of ChatGPT's', and they are encouraged to

ask the AI for further fine-tuning – more functional requirements that will help their product become more robust, as well as help them expand their personal vision and make it better. The process is showcased step-by-step in real-time, but a video is prepared in case of service interruption during the demonstration. Outages in the OpenAI services have been less frequent recently, but during the peak hours they still happen, and it is important to have a good backup ready. Creating these videos is like a snapshot in time – it also helps showcase the differences in versions and their progress, as well as the time it takes to answer queries, model differences (3.5 and 4.0), as well as answer consistency. Recording the results has created data points for further research.

A showcase on how the model can play its role is also being created but requires more testing. Researchers have emulated an entire software development team creating and finishing a project [2]. The idea is to showcase how the different roles in a software development process are represented by AI (e.g., "Imagine you are a business analyst working on project 8 - Airport Logistics Planner. What would your workflow be? How would you approach every one of your team members? Please describe your activities using a table."). Current focus is on the fact that users must highly customize their input and requirements gradually if they want to create and document a good and useful product.

**You**
Hello, please generate 10 ideas for a small software that has to be coded and documented. Put them in a table.

**ChatGPT**
Certainly! Here's a table with 10 software ideas along with a brief description for each:

| Idea Number | Software Title | Description |
| --- | --- | --- |
| 1 | Task Tracker | A simple task management tool with the ability to create, edit, and prioritize tasks with deadlines. |
| 2 | Expense Manager | An application to track and manage personal or business expenses, with features for budgeting and categorization. |
| 3 | Fitness Companion | A fitness app that allows users to log their workouts, track progress, and set fitness goals. |
| 4 | Recipe Organizer | An application for organizing and storing recipes, including features for meal planning and shopping lists. |
| 5 | Budget Analyzer | A financial tool that imports and analyzes bank statements to provide insights into spending habits and trends. |
| 6 | Language Learning Aid | A software to help users learn a new language, including vocabulary quizzes, pronunciation guides, and practice exercises. |
| 7 | Event Scheduler | An event management tool that enables users to schedule, coordinate, and manage events efficiently. |
| 8 | Habit Builder | An app focused on habit formation, allowing users to set goals, track daily habits, and receive reminders. |
| 9 | Mindfulness Assistant | A tool to support mindfulness and meditation practices, offering guided sessions and progress tracking. |
| 10 | Book Recommendation | A system that recommends books based on user preferences, reading history, and genre preferences. |

These ideas cover a range of categories, from productivity and health to entertainment and personal development. Choose the one that aligns best with your interests and skills!

**You**
Assume I am currently a student in Software engineering. I already know how to work with databases, how to write code and interpret algorithms. I am currently brainstorming for ideas to write a course project on a topic of my choosing. I would like you to generate 10 ideas for a small software project that has to be coded and documented. Put them in a table. I want the ideas to be in a field of my liking - I like computer games, planes, trains, I hate cars and busses, as well as logistics. I also like europe as a continent. Take my interests and the things I do not like into account when generating the ideas.

**ChatGPT**
Certainly! Considering your interests in computer games, planes, trains, and a dislike for cars and buses, as well as your fondness for Europe, here are 10 software project ideas tailored to your preferences:

| Idea Number | Project Title | Description |
| --- | --- | --- |
| 1 | Flight Simulator Companion | Develop a companion app for popular flight simulator games, offering in-flight checklists, weather updates, and navigation aids. |
| 2 | Train Timetable Tracker | Create a software tool to track and display train timetables across Europe, allowing users to plan and optimize their train travel. |
| 3 | EuroQuest: Geography Game | Design an educational game that quizzes users on European geography, landmarks, and capitals to make learning fun and interactive. |
| 4 | Virtual European Tour | Develop a virtual tour application that takes users on a journey through iconic European cities, showcasing historical landmarks and cultural hotspots. |
| 5 | Air Traffic Control Sim | Build an air traffic control simulation game where players manage air traffic at major European airports, emphasizing realistic scenarios and challenges. |
| 6 | Railway Tycoon Game | Create a strategy game where players can build and manage their own railway networks, dealing with challenges such as scheduling, maintenance, and competition. |
| 7 | European History Timeline | Develop an interactive timeline app that highlights key events in European history, providing users with a visually engaging way to explore historical periods. |
| 8 | Airport Logistics Planner | Design a logistics planning tool for airport operations, helping manage the flow of goods and resources efficiently within European airports. |
| 9 | Gaming Community Hub | Build a centralized platform for gamers interested in European-themed games, providing forums, reviews, and a curated list of games related to the region. |
| 10 | Rail Adventure Game | Create an adventure game set on a fictional train journey through Europe, incorporating puzzles, storytelling, and historical elements for an immersive experience. |

Feel free to choose the one that resonates most with you or sparks your interest!

Fig. 1. Comparison between statements for idea generation, showcasing the influence of detailed information to the outcome.

The next step comes when the students have decided on their ideas for each team. At this point they have written most of the requirements, their vision, and they have mostly decided on what technological stack they will use. A showcase on how to create a relational DB that covers their functional requirements is done.

On Fig. 2 the process can be seen – starting with asking for the creation of several versions of the relations in a database and GPT must draw them with tables and connect them. When solutions are iterated several times, the first version of the database is required is completed. At this point SQL statements that will create the DB and its relations are required, specifying the type of software and SQL server they are going to run on. The demonstration currently uses MS SQL server through the Management Studio, but other clients are viable alternatives too – like PostgreSQL or MySQL. The demonstrations do not yet include experimenting with communicating with ChatGPT for NOSQL solutions, due to time constraints, but they are worth looking into.



Fig. 2. Aasking ChatGPT with the creation of a database for a small-scale project.



Fig. 3. Inclusion of AI in the software development process in the ASRS subject.

Another important step is showing how to "seed" the database with relevant data. First, a request that simply returns random raw data into several queries is created (see Fig. 3). Afterwards, if necessary, several key points and examples that the message should include in order to get relevant data are defined – for example random data in a special regular expression format is requested. Another option is to give GPT a short list with the example data, or the pattern to which it conforms. GPT can then work with it and iterate on it. This is followed by asking for the insert statements that are used to seed the finished database. There are cases where it has been noticed that there are differences in the format of the required data, as well as problems with relations – the foreign keys are not populated in the corresponding order, which leads to issues. Those issues are easily solvable for students that have passed their database-related subjects and further support the learning experience in working with such tools.

Students are shown how to ask for UI suggestions (mainly HTML and CSS), but as their chosen project tech stacks are different (games, desktop applications, web sites and mobile apps), the use of AI in the development part of the project is concluded here. It is important to note that there are interesting applications for creating unit tests that are mentioned. This part of the subject curriculum is still in development so at this stage students are only shown how to make simple unit tests for C# and how to base them on their functional requirements.

To put the finishing touch to their projects, hints are given on using the roleplaying ChatGPT to help them with their final presentations. A dialogue is shown following a similar structure to: "Hello, you are a business analyst. You are going to present a project on the topic of a flight fleet management software system. Your audience consists of people working in the sector. The project is made using the following technologies: C#, MSSQL, WebSocket, JavaScript and others. What would be your plan to present this project if you only have 15 minutes to talk?". Students are usually very impressed, as most of them are used to using AI tools in other ways.

The use of AI powered software in the different steps of the software development process that is used in the ASRS subject can be seen in Fig. 3. As the focus is on creating requirements, specification and software documentation, it is acceptable to use AI for any of the other steps, as students have already learnt how to complete those by themselves. A strong

argument can be made that AI powered chatbots will generate the entire documentation themselves in the near future, but this stage of development is not yet reached. The guiding philosophy is that if students are able to finish a task by themselves and understand it well enough, they should be taught how to automate it at a future point in time.

### B. Inclusion of AI Assistance in the Artificial Intelligence Subject

The use of assistive tools with AI in the subject "Artificial Intelligence" is presented in the second half of the semester. Traditionally the first half of the workshop tasks for the students consist of using Python to solve classic AI problems, such as eight queens and different crosswords, while using search and constraint satisfaction algorithms. After the introductory problems are taught, several more advanced topics and workshops are covered. Students are shown neural networks (NNs) - single layered and multilayered. They are tasked with creating several NNs on their own using the Keras library and MNIST data sets. The use of convoluted NNs in various fields such as in finances - predicting financial time series [4], gives the opportunity to teach students the interdisciplinary aspects of software engineering in general. After all, what use is software engineering if there is no field to attach software to? Giving working examples in actual businesses and practical applications gives confidence and generates ideas in our pupils.

In one of the following workshops, instructions on how to use open-source face recognition libraries in Python to detect faces on photos of famous people are included. There are training sets provided that are used to train and validate the model. In the researchers' experience, doing this workshop takes a significant amount of time – both the process to fulfill all the necessary requirements on each of their personal machines, and then follow the instructions and test the code they have written step by step are time-consuming. After the students have managed to go through all the steps and have created a working piece of software, they are encouraged to train and validate the model using their own set of photos as training sets.

Facial recognition, and image recognition in general, is not a new concept, but it is imperative students are taught how to apply it and use LLM tools on their own product to improve and test it. Students gain more thorough experience on how to modify their software products more efficiently by having access to the ChatGPT terminal, and thus to the entire collective training set using hundreds of billions of parameters that OpenAI have been using and improving.

A simple visual representation of the approach to teaching students in the AI subject is shown in Fig. 4.

The process of preparing additional workshops is constantly undergoing, and some have yet to be completed with students, as the subject is next semester in their curriculum. One workshop includes a comparison in image generation methods for different assisted tools, several of the most popular AI tools are presented - Stable Diffusion, Bing Images, Bard and Dall-E. A sufficiently difficult prompt using weight distribution for each of its parameters is used.



Fig. 4. Methodology for practical application of theoretical knowledge base in the education of students.

It is worth noting that Bard also hallucinated and tried to say it "drew" an image and responded with a google image search during one of the attempts to make it draw.

Fig. 5, 6, 7 and 8 display some of the collated results that were achieved using the currently most popular available AI-powered models. After multiple attempts, only Bing produced acceptable results, while Bard simply could not create any images itself, and returned googled images in almost all the attempts. The online version of DALL-E could not make images of better resolution and quality than those shown on Fig. 5. It can be said that testing web services in this way is not extensive enough, but when the services are unreliable themselves, produce different results and are in most cases limited, paid or of bad quality, they cannot be reliably used for educational purposes.



Fig. 5. Using Stable diffusion multiple times using a prompt.



Fig. 6. Dall-E web using a prompt: "please draw a young robot that has a white coat, make the image ultra realistic".



Fig. 7. Bard's response to drawing prompts.

Fig. 8. Bing Images (Dall-E based).

It should be noted that Fig. 6 showcases the result of a very specific prompt with weight distribution that is as follows: "breathtaking, mysterious, fantasy, magical, (female robot:1.5) (creating a machine with hands), surrealism, hyper-realistic, colors and shapes, highly detailed, realism pushed to extreme, fine texture, 8k, ultra-detailed, (vivid swirling smoke, thick smoke:1.4), fluid, fire, cinematic, (intricate details:1.5), (vibrant colors:1.4), flash explosion, (colorful powder explosion), (billowing hair:1.5) "((magical enchantment on hands))" "((creates a machine from hands towards viewer))"". The LLM allows for incredibly detailed instructions that can be followed, such instructions are unavailable for the other tested tools currently. Fig. 5, 7 and 8 display the limitations of the most popular models. While spectacular results can also be achieved using them, it requires a bigger investment in prompting them, both in terms of time and money (as the paid model tokes are usually limited).

Some advantages and disadvantages are shown in Table I. The best results during the testing were achieved by StableDiffision. All the other services had either service disruptions (Bing, Bard), hallucinations (Bing, Dall-E), required tokens and accounts (Bing, Dall-E), or simply didn't do as requested (Bing, Bard). The multitude of disadvantages make them unfit to be used in higher education. They are worth looking into further if their future iterations are improved. The currently chosen AI for the future workshops is StableDiffision. Preliminary testing and scenarios based on the technology are being created and refined.

The methodology that has been adapted for an AI image generation workshop is shown on Fig. 9. It includes stable diffusion and showcases different models and their differences.

TABLE I. Comparison between the Different AI Powered Image Generation Tools

| AI | Advantages | Disadvantages | Notes |
|---|---|---|---|
| Bard | Free, on the web. Powered by Google. | Inconsistent behavior – sometimes says it cannot currently generate images; other times returns google image results. | Says it cannot generate images when prompted. After multiple attempts it starts to hallucinate that it can, in fact, generate images. The process itself is erroneous – it returns google image searches instead of generating images. |
| DALL-E | Free, on the web. Implemented in other software. | Low quality images, uses tokens, results are unsatisfactory. | Has bots in discord that work better than the software itself. |
| Bing | Free, uses Dall-e. | Only 5 questions. Has memory loss afterwards. Requires Microsoft Edge to work. Sometimes does not work – due to service disruption. | Says that it cannot generate images. Then generates images. |
| Stable Diffusion | Works offline. Can be downloaded for free. There are various community resources and models available. Active open-source project. | Requires hardware to run. Pre-trained models are at least several gigabytes each. Takes a long time to teach a model if hardware is not powerful enough. | While setting up the software can be difficult, the results are of great quality and require little manipulation before they can be used in production. |



Fig. 9. Methodology for showcasing AI image generation to software engineering students.

Stable diffusion has an advantage for the current purposes - meaning it can be locally installed for free by anyone, and it doesn't require services or subscriptions. It does have hardware requirements, but they can be satisfied relatively easily for lower image resolutions. There is an enormous number of pre-trained models that are available, as well as learning resources from its community. Several different approaches to combining or altering images after their creation are discussed and prepared for presentation.

### C. Introduction to Programming

For many students, the subject is the first time they work with algorithms and needing to understand what they are

doing. As it is taught during the first semester, the skill level difference is enormous. Preparing C++ tasks that are relevant and not demotivating is always a challenge everywhere in education [6, 7], some professors have very interesting solutions ways to solve the problem using games [8, 9]. The current established approach is to have automated tests for each small task we give the students, but it is still undecided on whether everyone should be introduced to testing at such an early semester. The benefits of test-driven development (TDD) are unquestionable [10, 11], but it requires a paradigm shift that is based on pre-existing knowledge in programming simple tasks, that many students lack during the first semester. Overloading fledgling software engineers with information is not the goal of the subject. The current idea is to use GPT to automate test creation process for the tasks that are given to the students for each workshop. The unit tests provide a good opportunity to limit test the code of students, but a reasonable way to enable them to access and understand them (as just giving a repository link is difficult for the average first-year student to understand) is not yet found. There are existing frameworks that provide a similar service [12], however they are not fit for this specific purpose.

### D. Computer Graphics

There have been discussions on how exactly to integrate tools such as Stable Diffusion and creating applications using technology based on API that works with an AI that is locally deployed, but the complexity of the task is daunting for students at the average skill level during the fifth semester. What has been learned is that there needs to be an in-depth course that helps students familiarize themselves with the usage of image generating software on a deeper level. Software engineers should be able to not only create images – the way they are taught during the Artificial intelligence course, but also understand AI image generating software and how to write code for applications that use the tool themselves. They need to be able to alter the image generating algorithms of the AI models, and create models fit to their projects' needs.

### E. Results and Observations

While the Dunning-Kruger effect has an impact on student self-assessment, research has shown that the differences between actual skill (and hence skill improvement) and self-evaluation are within acceptable levels [3]. It can be therefore proposed that students' self-assessment of their own skill improvement is significant enough to warrant noting a positive improvement due to the inclusion of AI in their curriculum subjects. Offloading the burden of menial tasks that can be easily automated leaves ample time for the actual improvement of their own professional interests and creating software products of value. One of the problems that have arisen is how professors evaluate whether the work done has led to improvement in a students' general skills, or just their skill in communicating with AI. That, however, is a complex enough problem to warrant research all by itself.

After questioning the students using a survey, as shown in Fig 10, it has been noticed that there is a good answer distribution among all the questions. It can be concluded that while it is useful, AI powered software needs improvement, and students using it need more experience to use it.



Fig. 10. Methodology for showcasing AI image generation to software engineering students.

What can be said with certainty is that students with the same university background have achieved success in projects that are more complex, more complete, and more in depth than those from previous years. According to the interviews with their teams during their project presentations, AI has certainly played a significant role in "doing the heavy lifting" for most of the teams. The most noticeable improvement is in the lower achieving students – given such powerful tools, they can elevate their skill level to a degree, sufficient for them to be motivated enough to be able to complete a reasonably sized software project. Compared to previous semesters in the same subject - students on the low skill spectrum could not finish as much work in the same amount of time. Elevating the results less skilled learners can achieve is motivating for professors as well. This gives confidence, experience, and motivation to the students, which in turn helps with student retention and an improvement to the system in higher education. Lifting the floor is beneficial to improving educational levels [5].

### III. CONCLUSION

Software engineering students need more exposure to AI tools in higher education. We need to have specialized subjects to help them familiarize themselves with this type of technology. It is imperative that they are not placed in a position where they have to "make do" but have the opportunity to learn how to leverage AI in their projects.

It is important to note that all LLMs often hallucinate, and because of that they are not reliable for checking facts and truth. At this point of their evolution, it is difficult to use them in the same way search engines are trustworthy data sources are used. What they are incredible at though, is giving ideas, workflows, helping with automation and doing the heavy lifting when beginning and polishing projects. This is what we are trying to teach our students, that they must be aware of the strengths and weaknesses of tools that apply AI.

With any emerging nascent technology, there is more testing required in order to make it work better. Familiarizing students with the correct approach to using AI assistive tools - focusing on their strengths and being aware of their weaknesses, is imperative for them to not be in a disadvantageous position in their future careers. As AI grows, everyone not aware of how to take advantage of it will quickly become a less valuable employee for the business. Starting too early is detrimental to brain development and problem-solving skills, but starting late is detrimental to career opportunities and life quality. Communication skills will become even more

important in the future but acquiring them will require more effort with the further digitization of communication and the lessening of social relationships people experience.

## REFERENCES

[1]  Oppenlaender, Jonas. "Prompt engineering for text-based generative art." arXiv preprint arXiv:2204.13988 (2022).

[2]  Qian, Chen, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. "Communicative agents for software development." arXiv preprint arXiv:2307.07924 (2023).

[3]  Gignac, Gilles E., and Marcin Zajenkowski. "The Dunning-Kruger effect is (mostly) a statistical artefact: Valid approaches to testing the hypothesis with individual differences data." Intelligence 80 (2020): 101449.

[4]  Markova, Maya. "Convolutional neural networks for forex time series forecasting." In AIP Conference Proceedings, vol. 2459, no. 1. AIP Publishing, 2022.

[5]  Crouch, Luis, and Caine Rolleston. "Raising the floor on learning levels: Equitable improvement starts with the tail." RISE Insights 2016 (2017): 1-14.

[6]  Aung, Shune Lae, Nem Khan Dim, Soe Mya Mya Aye, Nobuo Funabiki, and Htoo Htoo Sandi Kyaw. "Investigation of Value Trace Problem for C++ Programming Self-study of Novice Students." International Journal of Information and Education Technology 12, no. 7 (2022): 631-636.

[7]  Alzahrani, Nabeel, Frank Vahid, Alex Edgcomb, Kevin Nguyen, and Roman Lysecky. "Python Versus C++ An Analysis of Student Struggle on Small Coding Exercises in Introductory Programming Courses." In Proceedings of the 49th ACM Technical Symposium on Computer Science Education, pp. 86-91. 2018.

[8]  Ariffin, Mazeyanti Mohd, Nurshazlyn Mohd Aszemi, and Mohammad Syazran Mazlan. "CodeToProtect©: C++ programming language video game for teaching higher education learners." In Journal of Physics: Conference Series, vol. 1874, no. 1, p. 012064. IOP Publishing, 2021.

[9]  Agapito, Jenilyn L., Joshua C. Martinez, and J. D. Casano. "Xiphias: A competitive classroom control system to facilitate the gamification of academic evaluation of novice C++ programmers." In Proceedings of International Symposium on Computing for Education, ISCE, vol. 14, pp. 9-15. 2014.

[10] Beck, Kent. Test driven development: By example. Addison-Wesley Professional, 2022.

[11] Langr, Jeff. "Modern C++ Programming with Test-driven Development: Code Better, Sleep Better." Modern C++ Programming with Test-Driven Development (2013): 1-368.

[12] Markoska, Ramona. "Managing ICT solutions for training and evaluation of C++ programming skills in e-learning ecosystem." New Trends and Issues Proceedings on Humanities and Social Sciences 6, no. 7 (2019): 33-41.

# A Novel Approach to Data Clustering based on Self-Adaptive Bacteria Foraging Optimization

Tanmoy Singha[1], Rudra Sankar Dhar[2], Joydeep Dutta[3], Arindam Biswas[4]

Department of Electronic and Communication Engineering, National Institute of Technology, Aizawl, Mizoram, India[1, 2]
Department of Computer Science & Engineering, Siliguri Institute of Technology, Sukna, West Bengal[3]
School of Mines and Metallurgy, KNU, Asansol, West Bengal, India[4]

*Abstract*—Data clustering reduces the number of data objects by grouping similar data objects together. In this process, data are divided into valuable groups (clusters) or expressive without at all previous information. This manuscript represents a different clustering algorithm based on the technique of the adaptive strategy algorithm known as Self-Adaptive Bacterial Foraging Optimization (SABFO). It is a streamlining strategy for bunching issues where a cluster of bacteria forages to converge to definite locations as ultimate group communities by limiting the fitness function. The superiority of this method is assessed on numerous famous benchmark data sets. In this paper, the authors have compared the projected technique with some well-known advanced clustering approaches: the k-means algorithm, the Particle Swarm optimization algorithm, and the Fitness-Based Adaptive Differential Evolution (FBADE) Scheme. An experimental finding demonstrates the usefulness of the projected algorithm as a clustering method that can operate on data sets with different densities, and cluster sizes.

*Keywords—Data clustering; Self-Adaptive Bacterial Foraging Optimization (SABFO); Particle Swarm Optimization (PSO); FBADE scheme; the k-means algorithm and the classical BFO*

## I. INTRODUCTION

A method of analyzing and clustering unlabelled datasets is known as unsupervised machine learning. It is possible to find out the unknown patterns or data groupings using these algorithms without the involvement of human action. In the domain of cross-selling strategies, client segmentation image recognition, etc. The authors can employ the above strategies to find underlying patterns. It also enables us to discover similarities and differences in information. In short, a type of machine learning called unsupervised learning involves training models using unlabelled datasets and allowing them to act upon them without supervision. As opposed to supervised methods, clustering is an unsupervised method that works with datasets that do not have outcomes (target) variables or information about associations among explanations. The clustering algorithm is the key to data analysis and identifying groups (natural clusters). As a result of this process, similar data points are identified and grouped. Clusters make it easier to characterize the attributes of distinct entities. Users can then shift data and analyze certain categories as a result of this. Clustering allows organizations to address distinct client segments based on their attributes and similarities. This aids in profit maximization. If the dataset has too many variables, it can aid in dimensionality reduction. Irrelevant clusters can be detected and deleted from the dataset more easily.

For more than two decades, clustering has taken particular interest among scientists and researchers to apply in versatile domains. Researchers are continuously trying to develop better approaches to clustering.

In order to build knowledge-driven decisions, data mining provides upcoming behaviors that businesses can predict [1]. A clustering technique deals with discovering a structure in an unlabelled collection of data through unsupervised learning [2]. A data cluster is used in order to divide the enormous number of objects into smaller groups, so the objects with similar characteristics are clustered together, and the ones with dissimilar characteristics are in different groups [3]. A cluster is considered too many when it exceeds three, clustering becomes NP-complete, and it is, therefore, challenging to develop a well-organized clustering technique [4]. The clustering problem can be useful for segmenting images [5], clustering documents [6], predicting diseases [7], wireless-related sensors networks [8] analyzing common networks [9], identifying the traffic in the network [10], retrieving information [11], and marketing [12].Partitional clustering is being used in a numbers of real-life applications. It is an algorithm for dividing data objects into small groups based on defined criteria called the distance between them. Data samples are dispersed from one group to another iteratively according to the number of groups determined prior to implementation. To begin, it makes a set of partitions based on a definite measure. Data samples are divided into corresponding groups according to their centres [13]. K-means is the utmost widespread and fashionable algorithm among the partitional clustering algorithms as it is more practical and effective when handling heavy amounts of data. This algorithm's main drawbacks are its sensitivity to the initial cluster centres, inability to find global minima, and convergence to the local optima. Research on the improved BFO algorithm found that the algorithm has some limitations, including a fixed chemotactic step size and feeble bacterial connections.

As a result of the static chemotactic step size, it is problematic to strike the right balance among exploration and exploitation. Secondly, the feeble connection between bacteria shows a poor random city in chemotaxis. As a result of these two drawbacks, the bacteria community will search for a compound multimodal solution set at a local level rather than at a global level when compared to a global convergence.

An approach using self-adaptive BFO (SABFO) is proposed in this paper as one of the novelty of this article. It

improves the classical algorithm hypothetically in two ways. By leveraging bacterial search state features, the self-adaptive swimming process overcomes the traditional drawback caused by fixed step sizes. This paper extracts and calculates three vital qualities of the bacterial search state, namely the variety of the population, the number of iterations, and the mean fitness function.

The aim of this manuscript is to propose a new approach to clustering optimization technique, namely SABFO, which will provide performance optimization as the source of data grouping. A new perspective for solving NP-hard clustering problems is provided by Bacterial Foraging Clustering, a global optimization-based technique rather than high speed local search. At the same time, it's a new version of the Bacterial Foraging Optimization technique. In this proposed algorithm there is no need to select the centroid or center required to be chosen in the primary steps.

Further, the proposed algorithm aims to overcome two drawbacks of conventional algorithms:

*1)* The projected clustering technique achieved a high degree of accuracy as compared with other algorithms.

*2)* High-dimensional data can be processed resourcefully with the proposed algorithm.

The rest of this manuscript is planned in a subsequent way. In Section II, the literature review is introduced. Section III illustrates the preliminary knowledge of the BFO algorithm and optimization-based clustering. Section IV illustrates fully the entire method of the recommended SABFO-Clustering algorithm. In Section V, the authors present the numerical illustration for the datasets used in this paper. Sections VI and VII summarizes the results and discussion, respectively.

Finally, Section VIII shows the manuscript's conclusion and future work.

## II. LITERATURE REVIEW

The foraging performance of 18 Escherichia coli in the human being intestinal tract was studied, Passino [14] proposed a bacterium foraging optimization (BFO) algorithm in 2002 for optimizing 17 problems. Despite the 19 BFO algorithm's superiority over several other algorithms, 20 its convergence speed 21 and global search capability need to be improved, because it is extremely simple to fall into local optimal outcomes and convergence is slow.

Different clustering algorithms in recent years have been projected such as segmentation, density based hierarchical, grid-based, and model-based. By using partitioning, the authors can create partitions based on a number of criteria. The pattern belongs to only one cluster when using hard Partitional clustering. Clustering by fuzzy rules extends this notion by allowing patterns to belong to more than one cluster.

During the last few years, the BFO algorithm has often been combined 49 with other algorithms in various fields, Ofosu et al. In 50 [13], the proportional integral and derivative controllers 54 were unable to overcome the difficulties encountered in obtaining optimal PI gains 51 for fuzzy-PI controllers.

An optimal allocation model based on risk has been proposed by Xiong and et al. [15] and a multi-objective optimization57 problem has been solved by merging gradient particle, 58 swarm optimization with bacterial optimization reduces the risks associated with distributed 60 generation and facilitates the advancement of and implementation of distributed generation.

TABLE I.    LITERATURE REVIEW

| Author | Year | Technique introduced | Results |
|---|---|---|---|
| Tripathy, M and et.al [17] | 2006 | Enhanced Bacteria Foraging Optimization | Retained least cost |
| Li, M.S and et.al [18] | 2007 | Bacteria Foraging Algorithm varying population | Quorum sense, proliferation |
| Biswas, A and et.al [19] | 2007 | Genetic algorithm | Global optimization |
| Korani, W and et.al [20] | 2008 | PSO | Proportional – Integral – Derivative controller tuning |
| Dasgupta, S. and et.al [21] | 2009 | Micro Bacteria Foraging Optimization | Smaller population |
| Chen, H and et.al [22] | 2009 | Cooperative Bacteria Foraging Optimization | Explicit decomposition of search space |
| Dasgupta, S and et.al [23] | 2009 | Adaptive Bacteria Foraging Optimization | Varying chemotactic steps |
| Chen, H and et.al [24] | 2010 | Multi colony BFO | Several colonies |
| Kim, D.H and et.al [25] | 2011 | Genetic algorithm | Proportional – Integral – Derivative controller tuning |
| Gollapudi, S.V.R.S and et.al [26] | 2011 | PSO | Resonant frequency of rectangular micro strip antenna |
| Okaeme, N.A and et.al [27] | 2013 | Genetic algorithm | Automated investigational control design |
| Abd-Elazim, S.M and et.al [28] | 2013 | PSO | Power system stabilizers illustration |
| Mandeep Kaur and et.al [29] | 2018 | MOBFOA | Comparative study with other algorithm |
| Lv, X and et.al [30] | 2018 | IBFO | Machine Learning Framework |
| Huang Chen and et.al [31] | 2020 | SCBFO | Demonstration of CEC 2015 benchmark test set |
| Yufang Dan and et.al [32] | 2021 | BFO | Dynamic, multi-objective optimization, and complicated constrained optimization |
| Bo Yang and et.al [33] | 2022 | Discrete BFO | Unveiling global communities in networks |
| Sandeep Gogula and et.al [34] | 2023 | BFO | Size of the DGs, losses in active and reactive power flow |

Guo and Zhou [16] employed the trapezoid quadrature formula 65 in combination with the 64 BFO techniques to compute integrals since 64 integrable functions have many primitive functions that aren't elementary. Table I represents the tabular form of literature review with year, Technique and results used by the authors.

### III. PRELIMINARIES

#### A. BFO based Clustering

The process of clustering is a data mining method that involves classifying objects without any prior knowledge (clusters).It is possible to formalize the clustering problem as follows, given a sample data set $X=(x\_(1,) \ x\_(2 ,)……x\_(n )$ ),It is possible to formalize the clustering problem as follows, given a sample data set ,determine a partition of the objects into K clusters which satisfies:

$$\cup_{i=1}^{k} C_i = X; \qquad (1)$$

$$C_i \cap C_j \begin{cases} \emptyset \ , i,j = 1,2,………k; \ \ i \neq j \end{cases} \qquad (2)$$

$$C_i \neq \theta \quad i = 1,2,………k$$

In the mathematical point of view, cluster $C_i$ can be obtained by:

$$\begin{cases} |C_I = \{x_j | \|x_j - z_i\| \ \leq \ \|x_j - z_p\| \ , x_j \ \in X\} \\ q \neq 1, q = 1,2 ………….k \\ z_i = \frac{1}{|C_i|} \sum_{C_{j} \in C_i} x_j \ , i = 1,2,……….K \end{cases} \qquad (3)$$

Where,$\|.\|$ Signifies the length between of any two data points in the trial set, and $z_i$ =the centre of cluster $C_i$ .

#### B. The Classical BFO Algorithm

BFO algorithm is enlivened with a movement known as "chemotaxis" showed by bacterial foraging ways of behaving. Motile bacteria like E. coli and salmonella impel themselves by the turn of the flagella. An organic entity swimming or running forward is caused by the flagella turning counterclockwise, while a bacterium that makes a clockwise pivot tumble about with haphazard motion and swims once again. The bacterium can look for nutrients in any direction by switching among "swim" and "tumble" motions. The bacterium begins to swim more often as it gets closer to a nutritional gradient. Bacteria move away from some nourishment to search for further, resulting in tumbling, hence direction changes. Chemotaxis, in its simplest form, involves bacteria swimming and tumbling to reach advanced concentrations of food.

#### C. Bacterial Foraging Optimization

The three main mechanisms that make up the classical BFO system are chemotaxis, reproduction, and elimination-dispersal. Here are quick summaries of each of these processes:

*1) The basic chemotaxis:* Chemotaxis for bacteria is the course of the accumulation to nutrient-enriched regions. Bacterial movement designs incorporate both tumbling and swimming. Flips are the unit step lengths that bacteria take when moving in any direction. The extent to which adjustments are necessary determines whether the new position is more attractive than the opposite position. Then, the bacterium will keep up shifting in a more excited propensity for a couple of steps till the limit for variation is not any more shut. The enhanced method will be

$$Q_i \ (j + 1, k, l) = \ Q_i \ (j, k, l) + \frac{\Delta_i}{\sqrt{\Delta^T (i)\Delta_i}} \ C(i)n \qquad (4)$$

Here, $Q_i \ (j + 1, k, l)$ symbolize the $i$th bacterium at the $j$th chemotactic, kth denotes the reproductive, and lth represented the elimination dispersal steps. $C(i)n$ is denoted as the trend step length of bacteria $i$ in a random direction and $\Delta$ lies between −1 and 1 as a random vector.

*2) Swarming:* The chemotactic behavior of bacteria is not limited to searching for food individually but also includes both gravitation and repulsion between them in the foraging process. A bacteria's attractive information causes it to move to the center of the population, thus fetching the bacteria closer together. Yet, the bacteria's repulsion information keeps them at a distance from each other at the same time.

*3) Reproduction:* Eventually, bacteria with weak feeding abilities will be removed, while bacteria with robust feeding capabilities evolve to breed offspring to continue the population size. This process follows the usual method of survival of the fittest. The authors proposed a reproduction operation based on simulating this phenomenon. A chemotaxis operator performed by S/2 bacteria eliminated bacteria with poor fitness and let those with higher fitness self-replicate in S-sized populations.

A completed reproduction operation ensures that the offspring inherits the superior characteristics of the parents, and it also results in the protection of the good individuals and the acceleration of the progress towards an optimal global outcome. Fig. 1 represents the basic structure of the Bacteria Foraging Optimization algorithm.



Fig. 1. The basic organization of the Bacteria Foraging Optimization algorithm.

*4) Elimination:* It is important not to rule out the possibility that unexpected conditions might cause bacteria to die or migrate to a new location during bacterial foraging. It has been proposed to model this phenomenon by simulating elimination-dispersal operations.

## IV. PROPOSED SELF-ADAPTIVE BFO (SABFO)

Chemotaxis is a crucial tool in exploring and exploitation of the BFO algorithm during the search for the optimization space. The consequence of chemotaxis depends on the size of the steps and the direction the swimmer flips.

Two improvements are proposed in this paper to get better performance of the BFO algorithm, including extracting and calculating the features of the search state and increasing bacteria's communication. The SABFO algorithm provides a novel BFO algorithm to design dynamic self-adaptive swimming and flipping motions for bacterial cells with these two improvements.

It seems that the method of calculating swimming step size depends on the single swimming step size, C(t), as well as the quantity of chemotaxis, n, as indicated by the above researcher. An algorithm's performance can be adjusted to the ebb and flow state of search based on the size of the swimming steps. At the point when the pursuit state is in the beginning phase, the calculation needs the investigation capacity for worldwide pursuit; then, at that point, in the later stage, the abuse capacity is expected for nearby turn of events.

For various optimization issues, the difference in the BFO search state is likewise unique. In the meantime, on the grounds that the chemotaxis method is nonlinear and it is very complex to such an extent that the progress from the worldwide investigation to the nearby double-dealing can't be essentially depicted and separated by the way of logical conditions.

The proposed paper separates the three components of the BFO in every search state so that the authors can better understand the unique change of the BFO algorithm to the suitable chemotaxis swimming, including iteration, population diversity, and mean fitness. Fig. 2 represents the bacterial population position of BFO Algorithm.



Fig. 2. The bacterial population position of BFO algorithm.

The population diversity of bacteria refers to where it is dispersed. A wider range of bacteria dispersing will increase

population diversity, and vice versa. As part of the chemotaxis process of BFO, this manuscript measures the bacterial colony's population diversity.

$$div(t) = \frac{1}{D*S} \, g\sqrt{\sum_{i=1}^{S}\left(\frac{Q_i(j,k,l) - \overline{Q_t(J,k,l)}}{|L|}\right)} \qquad (5)$$

Where div (t) is the range between [0, 1], L denotes the solution space's elongated radius. This method measures the distance between the center of each bacterium and the solution space, regardless of the amount or dimension of the bacteria.

The iteration of the BFO algorithm is communicated through a boundary T, which is characterized as an articulation in the reach of [0,1] in equation 6, where t and $T_{max}$ address the record of the current chemotaxis and the most extreme emphasis, separately. Hence, the meaning of parameter T is for the most part suitable to various algorithms regardless of how the parameters, the aspect, with the arrangement space which are programmed into the algorithms:

$$T(t) = \frac{t}{T_{max}} \qquad (6)$$

The modification in the mean fitness function is in two chemotaxis processes, where dJ is primarily calculated as one of the essential values for examining the BFO algorithm. To provide a common explanation, the difference in the mean fitness dJ is characterized in the per-unit structure inside [−1,1] as follows:

$$dJ(t) = \frac{J(t)-J(t-1)}{J_{Max}-J_{Min}} \qquad (7)$$

Where, $J_{Max}$ and $J_{Min}$ denotes the maxima and minima of the fitness function, respectively. Fig. 3 shows the flowchart of the proposed algorithm SABFO.



Fig. 3. The flowchart of the proposed algorithm.

So, in this proposed manuscript, the 03 most important variables are the population diversity, the number of iterations and the bacteria's mean fitness value is used as input, which is intended to examine the pursuit status of the algorithm in these manuscripts. As the chemotaxis will be changed as per accompanying with the swim processes are as follows:

$$\begin{cases} n(t+1) = n(t) + dn(t) \\ C(t+1) = C(t)gC_M(t) \end{cases} \quad (8)$$

Where, the 02 output variables are the bacterial swimming movement increases $dn(t)$ [0.01, 1].as well as the bacterial swimming step multiple is C(t) (0, 1).

Another important operation of the BFO algorithm is flipping the bacteria. The direction of chemotaxis is determined by the extremum of each bacterium when swimming. Despite its benefits to the randomness of the search, this method results in a slower search because of blocked information among the bacteria. Therefore, BFO algorithms with suffer from the drawback of tumbling into the local optimum.

In order to resolve the issue, mentioning individuals' data exchange information strategy is used in BFO. So, the BFO algorithm is updated and the flipping variable is given by

$$\Delta_{t+1}(i) = wg\Delta_t(i) + C_1 R_1(P_{local} - P_{N_c}) + C_2 R_2(P_{global} - P_{N_c}) \quad (9)$$

By adjusting the co-efficient, the chemotaxis process is used in the above-mentioned equation.

Where w, denotes the chemotaxis inertia of the bacteria at a specific distance.

During bacterial chemotaxis, C1 signifies the amount at which each bacterium travels toward its distinct optimal value $P_{local}$, whereas C2 records the global optimum value $P_{global}$ for all bacteria. For improving the randomness of bacterial flipping and enhancing search ability, R1 and R2 are random the values between 0 and 1. The flowchart of the algorithm is shown in the Fig. 3.

## V. NUMERICAL ILLUSTRATION

In order to compare this proposed approach to Self-Adaptive BFO, five real-life data are Iris [35], Glass [35], breast cancer [35], Wine [35] and Vowel Dataset [35] were used in this proposed paper.

Real- Life Data Sets

- Iris Data: It comprises of three distinct types of iris blossom.

- Glass: The information was examined from six dissimilar kind of glass.

- Breast cancer: It comprises of 9 applicable highlights.

- Wine Data set: It is the outcome of a chemical analysis of wine. This analysis is resolved with the quantities of 13 constituents shown in every one of the three kinds of wine.

- Vowel Dataset: It comprises of 871 Indian Telugu vowel sounds.

The authors have used the following real-life datasets in this proposed paper. Table II represents the data set used in the manuscript.

TABLE II. THE DATASETS USED

| Real-Life Dataset | n | D | K |
|---|---|---|---|
| Iris plants [35] | 150 | 4 | 3 |
| Glass[35] | 214 | 9 | 6 |
| Wisconsin breast Cancer data set[35] | 683 | 9 | 2 |
| Wine[35] | 178 | 13 | 3 |
| Vowel Dataset[35] | 871 | 3 | 6 |

Where, n represents number of data points. D represents number of features and K represents no. of Cluster.

Similarly Table III represents the value of parameter used in the article.

TABLE III. VALUE OF PARAMETER

| Algorithm Name | Parameter Name | Value |
|---|---|---|
| FBADE | Population size | 10*dim |
| | Crossover | 0.9 |
| | Mutation | 0.8 |
| | $K_{max}$ | 20 |
| | $K_{min}$ | 2 |
| K Mean | Population size | 50 |
| | $\mu_c$ | 8 |
| | $\mu_m$ | 0.001 |
| | $K_{max}$ | 20 |
| | $K_{min}$ | 2 |
| PSO | Population size | 100 |
| | Inertia Weight | 0.72 |
| | $C_1, C_2$ | 1.494 |
| | $P_{initial}$ | 0.75 |
| | $K_{maximum}$ | 20 |
| | $K_{minimum}$ | 2 |
| SABFO | Population size (S) | 50 |
| | $N_C$ (No. of Chemotactic Steps) | 100 |
| | $N_S$ (Length of One swim) | 4 |
| | $N_{re}$(No. of reproduction steps) | 4 |
| | $N_{ed}$ (No. of elimination dispersal events) | 2 |
| | $P_{ed}$ (Probability of elimination dispersal events) | 0.25 |

## VI. RESULT ANALYSIS

Three performance metrics have been used to compare the SABFO algorithm with other evolutionary algorithms as state-of-the-art clustering techniques:

*1)* Performance metrics in the CS and DB domains as well as the number of misclassified items for each dataset;

*2)* Finding the optimal number of clusters; and

*3)* Computing time.

It is first necessary to measure the fair time of stochastic algorithms like PSO, FBADE, SABFO, and K Mean in order to compare their speed. Meanwhile the algorithms perform a dissimilar amount of work within their inner loops, as well as having different populations, the number of runs or generations cannot be used as a time measurement. Thus the authors choose to calculate computation time based on the number of fitness function evaluations (FEs) rather than the number of generations and iterations. When function complexity increases, counting the FEs is a reliable gauge of runtime complexity because it corresponds strongly with actual processor time.

Generally, two successive runs of four competing algorithms do not match because they are stochastic. As a consequence, the authors conducted 50 independent runs of each algorithm using different seeds. Based on the 40 runs, each result is expressed as a mean and standard deviation. As the various leveled agglomerative calculation utilized here, utilizes no developmental strategy, the amount of function evaluations isn't applicable to this technique. In this algorithm, the authors use the Ward updating equation to efficiently calculate cluster distances given the number of clusters for each problem.

Depending on the clustering validity measure used, any of the four evolutionary clustering algorithms will perform well. A CS measure-based fitness function is used in one set of experiments, while a DB measure-based fitness function is used in the other set of experiments. Four partitional clustering algorithms have been evaluated in terms of CS and DB calculation against the average-link metric based hierarchical method for each dataset.

This algorithm was run in Matlab 2010 under Windows 11 using an Intel Core i5 computer having 3.60 GHz speed and 8 GB of RAM.

The SABFO algorithm continues to offer superior clustering accuracy to each of the other three competitors as shown in Table IV. Tables IV and V represent the first four evolutionary algorithms (using the CS measure), mean classification error and standard deviation over nominal partitions were determined over 40 independent runs.

TABLE IV.    $10^6$ Function Evaluations (FES) with Cluster Strictness (CS)

| Name of the Dataset | Algorithm | Avg No. of clusters found | Value of CS calculated | Mean Intra cluster Distance | Mean Inter cluster Distance |
|---|---|---|---|---|---|
| BreastCancer | SABFO | **2.26±0.00** | **0.4623±0.033** | **4.2356±0.143** | **3.2489±0.138** |
|  | PSO | 2.13±0.0587 | 0.4878±0.009 | 4.7845±0.356 | 2.3521±0.021 |
|  | K Mean | 2.00±0.0079 | 0.5098±0.015 | 4.8879±0.904 | 2.3857±1.699 |
|  | FBADE | 2.06±0.0232 | 0.4854±0.359 | 4.5944±0.599 | 2.8977±1.345 |
|  | Classical BFO | 2.15±0.0261 | 0.8984±0.381 | 4.5644±0.546 | 3.0625±1.455 |
| Vowel | SABFO | **5.72±0.0641** | **0.9068±0.046** | **1399.96±0.692** | **2698.58±0.112** |
|  | PSO | 7.25±0.0183 | 1.1827±0.431 | 1482.51±3.973 | 1923.93±1.154 |
|  | K Mean | 5.05±0.0075 | 1.8978±0.897 | 1485.13±12.235 | 1921.38±0.742 |
|  | FBADE | 7.50±0.0569 | 1.0844±0.067 | 1493.72±10.833 | 2434.45±1.213 |
|  | Classical BFO | 6.66±0.0895 | 1.2335±0.048 | 1499.96±0.956 | 2698.58±0.112 |
| Glass | SABFO | **6.04±0.0139** | **0.3221±0.456** | **563.247±134.2** | **853.62±9.044** |
|  | PSO | 5.89±0.0093 | 0.7532±0.073 | 599.535±10.34 | 889.32±4.233 |
|  | K Mean | 5.82±0.0346 | 1.4743±0.236 | 594.673±30.62 | 869.93±1.789 |
|  | FBADE | 5.59±0.0754 | 0.6999±0.643 | 608.787±20.92 | 891.82±4.945 |
|  | Classical BFO | 5.89±0.0654 | 0.7506±0.725 | 598.852±166.3 | 890.89±8.250 |
| Iris | SABFO | **3.198±0.0382** | **0.6548±0.097** | **3.106±0.033** | **2.3941±0.027** |
|  | PSO | 2.23±0.0443 | 0.7361±0.671 | 3.6516±1.195 | 2.2104±0.773 |
|  | K Mean | 2.35±0.0985 | 0.7282±2.003 | 3.5673±2.792 | 2.5058±1.409 |
|  | FBADE | 2.50±0.0473 | 0.7633±0.039 | 3.9439±1.874 | 2.1158±1.089 |
|  | Classical BFO | 2.65±0.0752 | 0.7531±2.003 | 3.6689±1.562 | 2.2515±1.233 |
| Wine | SABFO | **3.19±0.0391** | **0.8989±0.032** | **4.041±0.002** | **3.1399±0.078** |
|  | PSO | 3.03±0.0253 | 1.7899±0.037 | 4.787±0.184 | 2.6113±1.637 |
|  | K Mean | 2.95±0.0112 | 1.5842±0.328 | 4.163±1.929 | 2.8058±1.365 |
|  | FBADE | 3.50±0.0143 | 1.7964±0.802 | 4.949±1.232 | 2.6118±1.384 |
|  | Classical BFO | 3.65±0.0562 | 1.6998±0.056 | 4.655±0.095 | 2.922±1.563 |

TABLE V.    MEAN CLASSIFICATION ERROR

| Dataset | Mean Classification Error | | | | |
|---|---|---|---|---|---|
| | SABFO | PSO | K Mean | FBADE | Classical BFO |
| Breast Cancer | 21.98±0.28 | 27.01±1.25 | 29.00±1.55 | 29.15±0.50 | 26.00±0.00 |
| Vowel | 413.88±3.08 | 451.58±5.98 | 471.69±6.89 | 474.72±4.25 | 496.00±0.00 |
| Glass | 91.51±0.19 | 102.1±0.68 | 98.21±0.08 | 105.36±0.54 | 111.00±0.00 |
| Iris | 2.35±0.00 | 4.15±0.0 | 5.00±0.00 | 3.96±0.00 | 4.00±0.00 |
| Wine | 36.52±0.0 | 98.4±1.09 | 100.24±1.05 | 114.50±1.53 | 134.00±0.00 |

TABLE VI.    DB VALUES AT THE PREDEFINED CUT-OFF VALUE WERE CALCULATED AFTER 50 INDEPENDENT RUNS, AND THE MEAN CLASSIFICATION ERROR

| Name of the Dataset | Name of the Algorithm | Mean no. of FE's required | DB Cutoff Value | Mean Intra cluster Distance | Mean Inter cluster Distance |
|---|---|---|---|---|---|
| Iris | SABFO | **504783.45**±12.65 | 0.8 | **3.9928±0.029** | **2.1029±0.842** |
| | PSO | 679084.75±16.57 | | 3.7852±1.842 | 1.7641±0.439 |
| | K Mean | 790865.90±10.21 | | 4.4587±3.782 | 1.9383±1.307 |
| | FBADE | 658796.3 | | 4.0393±1.5 | 1.6278±1.6 |
| Wine | SABFO | **464653.35**±5.50 | 6 | **4.8292±0.732** | **3.0219±0.069** |
| | PSO | 486885.85±2.85 | | 5.1472±0.472 | 2.1161±1.623 |
| | K Mean | 598743.35±8.09 | | 4.9383±1.722 | 2.9121±0.353 |
| | FBADE | 477869.95±8.12 | | 4.7531±2.043 | 2.8158±0.389 |
| Breast-Cancer | SABFO | 424732.30±8.93 | 0.9 | 5.4489±0.342 | 3.0234±0.683 |
| | PSO | 467854.60±10.12 | | 5.2885±0.552 | 2.0124±1.596 |
| | K Mean | 678874.90±7.82 | | 6.8832±0.733 | 2.1637±1.458 |
| | FBADE | **418765.55±1.23** | | **5.8684±0.467** | **1.9235±0.164** |
| Vowel | SABFO | **435743.05±2.65** | 3 | **1544.92±0.834** | **2081.31±0.679** |
| | PSO | 556865.00±4.26 | | 1652.58±2.341 | 1264.87±3.069 |
| | K Mean | 575854.65±1.29 | | 1582.55±7.332 | 1989.38±7.734 |
| | FBADE | 546859.60±2.05 | | 1608.22±5.866 | 1604.43±1.674 |
| Glass | SABFO | **506754.00**±12.27 | 2 | **132.757±15.8** | **13.46±2.54** |
| | PSO | 569787.95±10.83 | | 154.564±39.6 | 13.56±2.65 |
| | K Mean | 687678.75±10.97 | | 155.856±24.7 | 10.42±4.69 |
| | FBADE | 527585.35±7.50 | | 178.809±30.3 | 10.21±1.09 |

CS and DB index values were reduced by the SABFO within a minimum number of function evaluations in the majority of cases, as shown in Tables VI. According to Table VI, SABFO continues to provide superior clustering accuracy to the other three competitors. Entries of Statistically significant differences between SABFO and its competitors are evident in Table VI, only for breast cancer, FBADE yield a lower DB value than SABFO. Table VII shows the mean classification error and standard deviation of the different data set.

Table VIII represent the first four evolutionary algorithms (using the DB measure), mean classification error and standard deviation over nominal partitions were determined over 40 independent runs.

TABLE VII.    MEAN CLASSIFICATION ERROR AND STANDARD DEVIATION

| Dataset | Mean Classification Error | | | | |
|---------|------|-----|--------|-------|--------------|
| | SABFO | PSO | K Mean | FBADE | Classical BFO |
| Iris | **2.22±0.00** | 2.79±0.55 | 2.75±0.08 | 2.74±0.00 | 3.14±0.00 |
| Wine | **40.15±0.0** | 112.5±2.50 | 118.45±1.77 | 76.45±0.236 | 102.22±1.05 |
| Breast Cancer | **26.72±0.25** | 30.33±0.48 | 26.55±0.79 | 29.00±1.12 | 29.03±1.09 |
| Vowel | **416.37±7.50** | 437.00±3.72 | 476.58±3.59 | 478.62±2.69 | |
| Glass | **8.86±0.42** | 14.35±0.26 | 17.98±0.67 | 15.69±0.85 | |

TABLE VIII.    DB MEASURE-BASED FITNESS FUNCTIONS

| Name of the Dataset | Name of the Algorithm | Average Number of clusters found | Value of DB calculated | Mean Intra cluster Distance | Mean Inter cluster Distance |
|---------------------|-----------------------|----------------------------------|------------------------|-----------------------------|-----------------------------|
| Iris | SABFO | **3.48±0.0217** | **0.4644±0.029** | **3.1636±0.078** | **2.8389±0.678** |
| | PSO | 2.28±0.0598 | 0.6677±0.008 | 3.8536±0.122 | 2.2548±0.034 |
| | K-Mean | 2.32±0 | 0.7269±0.0 | 3.8428±0.076 | 2.1438±0.020 |
| | FBADE | 2.51±0.0089 | 0.5825±0.069 | 3.8879±0.089 | 2.0358±0.058 |
| | Classical BFO | 2.96±0.008 | 0.8674±0.00 | 3.8098±0.00 | 2.2857±0.00 |
| Wine | SABFO | **3.25±0.0931** | **3.0432±0.021** | **4.4212±0.096** | **3.1029±0.047** |
| | PSO | 3.05±0.0024 | 4.3432±0.232 | 4.8668±0.154 | 2.6113±1.635 |
| | K-Mean | 2.95±0.0173 | 5.3424±0.343 | 5.1312±1.342 | 2.7565±2.128 |
| | FBADE | 3.50±0.0143 | 3.3923±0.092 | 4.263±1.907 | 2.8158±1.786 |
| | Classical BFO | 2.99 | 5.7206±0.00 | 4.982±0.00 | 2.5009±0.00 |
| Breast Cancer | SABFO | 2.48±0.0653 | 0.5102±0.007 | 4.5564±0.024 | 3.1020±0.068 |
| | PSO | 2.50±0.0621 | 0.5754±0.073 | 4.9232±0.373 | 2.2684±0.063 |
| | K-Mean | 2.50±0.0352 | 0.6328±0.002 | 6.5541±0.433 | 1.8032±0.016 |
| | FBADE | **2.10±0.0081** | **0.5199±0.007** | **5.2234±0.042** | **2.0236±0.058** |
| | Classical BFO | 2 | 0.7634±0.00 | 5.0098±0.00 | 2.2817±0.00 |
| Vowel | SABFO | **5.75±0.0241** | **0.9224±0.334** | **1449.12±0.834** | **2289.85±0.163** |
| | PSO | 7.25±0.0562 | 1.2821±0.009 | 1500.57±3.748 | 1747.76±1.764 |
| | K-Mean | 5.05±0.0561 | 2.9482±0.028 | 1573.23±4.675 | 2271.89±1.222 |
| | FBADE | 7.50±0.0819 | 1.4488±0.075 | 1498.78±2.725 | 1962.31±0.993 |
| | Classical BFO | 6 | 3.0581±0.00 | 1493.98±0.00 | 2357.62±0.00 |
| Glass | SABFO | **6.05±0.0248** | **1.0092±0.083** | **501.757±4.3** | **893.46±3.32** |
| | PSO | 5.95±0.0193 | 1.5152±0.073 | 514.554±9.5 | 856.00±8.07 |
| | K-Mean | 5.85±0.0346 | 1.8371±0.034 | 518.903±2.9 | 852.32±5.43 |
| | FBADE | 5.60±0.0446 | 1.6673±0.004 | 514.849±3.4 | 862.21±2.53 |
| | Classical BFO | 6 | 1.8519±0.00 | 610.033±0.00 | 895.47±0.00 |

TABLE IX.    MEAN CLASSIFICATION ERROR

| Dataset | Mean Classification Error | | | | |
|---------|------|-----|--------|-------|--------------|
| | SABFO | PSO | K-Mean | FBADE | Classical BFO |
| Iris | **2.21±0.02** | 2.80±0.56 | 2.78±0.10 | 3.15±0.07 | 2.75±0.01 |
| Wine | **41.25±0.01** | 112.5±2.50 | 118.45±1.77 | 103.20±1.05 | 58.15±0.08 |
| Breast Cancer | **27.69±0.28** | 30.23±0.46 | 26.50±0.80 | 29.00±1.09 | 29.08±0.25 |
| Vowel | **417.39±6.99** | 435.00±3.75 | 473.46±3.57 | 472.65±2.76 | 486.65±3.26 |
| Glass | **8.82±0.42** | 14.56±0.28 | 17.98±0.67 | 15.70±0.89 | 17.52±0.68 |

Fig. 4.   The 3D plot of the unlabeled Iris data set.

The authors have applied various well-known advanced clustering approaches like the k-means algorithm, the Particle Swarm optimization algorithm, and the Fitness-Based Adaptive Differential Evolution (FBADE) Scheme on the 3D plot of Iris data (Fig. 4). The clustering results are as follows:



Fig. 5.   Clustering of iris data by SABFO.



Fig. 6.   Clustering of iris data by PSO.



Fig. 7.   Clustering of iris data by K-Mean.

Fig. 5 can classify the data set which contain overlapped clusters very efficiently and it also the ability to cluster data sets with high dimension as compared to Fig. 6, 7 and 8.



Fig. 8.   Clustering of iris data by FBADE.



Fig. 9.   The 1D plot of the unlabeled Wine data set.

Fig. 9 represents the 1-Dimensional plot of the unlabeled wine data set. The authors have also applied various well-known superior clustering approaches like the k-means algorithm, the Particle Swarm optimization algorithm, and the Fitness-Based Adaptive Differential Evolution (FBADE) Scheme on the 1D plot of Wine Data set. The clustering results are as follows:



Fig. 10. Clustering  of unlabeled Wine data set by SABFO.



Fig. 11. Clustering  of unlabeled Wine data set by PSO.

Fig. 12. Clustering of unlabeled Wine data set by K-Mean.



Fig. 13. Clustering of unlabeled Wine data set by FBADE.

Fig. 10 can classify the data set which contain overlapped clusters very efficiently and it also the ability to cluster data sets with high dimension as compared to Fig. 11, 12 and 13.

## VII. DISCUSSION

As can be seen in Table IV, the SABFO algorithm continues to offer clustering accuracy that is superior to that of the other three competitors. The first four evolutionary algorithms are shown in Tables IV and V. The mean classification error and standard deviation over nominal partitions were calculated after 40 independent runs using the CS measure.

CS and DB record values were decreased by the SABFO inside a base number of capability assessments in most of cases, as displayed in Tables VI. According to Table VI, SABFO keeps on giving better bunching exactness than the other three contenders. Passages of Genuinely tremendous contrasts among SABFO and its rivals are obvious in Table VI, just for bosom disease, FBADE yield a lower DB esteems than SABFO.

Clustering problems that have several data items, clusters, and overlapping cluster shapes have noticeable performance changes. The clustering accuracy of SABFO is consistently superior to that of its competitors in both Tables IV and V. The FBADE and SABFO methods have two clusters nearly same every time when it's run for the breast cancer dataset, despite having very similar final CS indices. Entries of Statistically significant differences between SABFO and its competitors are evident in Table VI, only for breast cancer, FBADE yield a lower DB value than SABFO.

The results of Tables V and IX indicate that the SABFO produces the fewest misclassified items after clustering. Although all five algorithms demonstrated convincing performance, there were misclassifications in each experiment based on the nominal classification, as expected. In this proposed evolutionary clustering algorithms, the authors found that the fitness values obtained were much better than those obtained from the insignificant classification, which represents that optimization cannot explain through misclassification. As a result, misclassification is caused by underlying expectations in the clustering fitness values (such as clusters' spherical shape), outliers in the dataset, and errors in data collection and nominal solutions. This is indeed not a negative result. Clustering solutions based on statistical criteria and minor classifications can be compared to reveal interesting data points and anomalies. Using a clustering algorithm to pre-analyze data in this way can be very useful.

According to Tables IV and VIII, both the CS and DB indices reached their cut-off values within a minimum number of FE's.

## VIII. CONCLUSION

A SABFO algorithm was proposed in this manuscript to address the fixed step size of the classical BFO algorithm as well as weak correlation among bacteria. Self-adaptive chemotaxis is an adaptation of the self-adaptive swimming technique depends on bacteria's state of search features, combined with enhancement of chemotaxis flipping based on exchange of information.

A comparison of the SABFO algorithm on 05 data sets was conducted by the PSO algorithm, the FBADE algorithm, the K-Mean algorithm and the classical BFO. It was found that SABFO's algorithm is accurate and effective at determining optimal solutions based on the validation results. The SABFO algorithm was also demonstrated to have better exploitation abilities in the future stages and having a more steady search performance.

In brief, the SABFO algorithm has a good stability between exploration and exploitation, which reduces the risks of local convergence. Further it can overwhelm the aforesaid two shortcomings of traditional BFOs.

Additionally, the SABFO algorithm is very stable and performs well when searching. Due to this, SABFO provides an efficient and novel way to accord with complex optimization issues.

In the above table, it appears that all four competitor algorithms terminated with similar accuracy for all the datasets. Based on the proposed algorithm, the CS and DB are found very lowest as per Table IV and VIII. In addition, SABFO successfully found the near-exact number of classes over consecutive iterations (three for iris and Wine data sets).For future researchers, there is a lot of scope to improve the proposed variants that may give much more excellent results on real-world optimization problems.

### REFERENCES

[1] P. P. Mohanty, S. K. Nayak, U. M. Mohapatra, and D. Mishra 2019 A survey on partitional clustering using single-objective metaheuristic approach. Int. J. Innovative Comput. Appl. vol.10 pp.207-226.

[2] A. Abraham, S. Das, and S. Roy 2008 Swarm intelligence algorithms for data clustering in Soft computing for knowledge discovery and data mining Springer pp.279-313.

[3]    T. Niknam and B. Amiri, 2010 An efficient hybrid approach based on PSO, ACO and k-means for cluster analysis Appl. Soft Comput. Vol.10 pp.183-197.

[4]    G. Sahoo 2017 A two-step artificial bee colony algorithm for clustering Neural Comput. and Appl. vol.28 pp.537-551.

[5]    A. Nithya, A. Appathurai, N. Venkatadri, D. Ramji, and C. A. Palagan 2020 Kidney disease detection and segmentation using artificial neural network and multi-kernel k-means clustering for ultrasound images Measurement vol.149.

[6]    A. Christy and G. M. Gandhi, 2020 Feature Selection and Clustering of Documents Using Random Feature Set Generation Technique in Advances in Data Science and Management: Springer pp.67-79.

[7]    S. Ramasamy and K. Nirmala, 2020 Disease prediction in data mining using association rule mining and keyword based clustering algorithm Int. J. Comput. Appl. vol.42 pp.1-8.

[8]    D. K. Kotary and S. J. Nanda, 2020 Distributed robust data clustering in wireless sensor networks using diffusion moth flame optimization Engineering Applications of Artificial Intelligence vol.87 First International Conference on Advances in Physical Sciences and Materials Journal of Physics: Conference Series 1706 (2020) 012163 IOP Publishing doi:10.1088/1742-6596/1706/1/01216310.

[9]    S. A. Curiskis, B. Drake, T. R. Osborn, and P. J. Kennedy, 2020 An evaluation of document clustering and topic modeling in two online social networks: Twitter and Reddit Inf. Process. Lett. & Manage. vol.57.

[10]   J. Yang, Y. Han, Y. Wang, B. Jiang, Z. Lv, and H. Song, 2020 Optimization of real-time traffic network assignment based on IoT data using DBN and clustering model in smart city Future Gener. Comput. Syst. vol.108 pp.976-986.

[11]   P. Bedi and S. Chawla, 2010 Agent based information retrieval system using information scent Int. J. Artif. Intell. vol.3 pp.20-238.

[12]   B. Yue, 2020 Topological Data Analysis of Two Cases: Text Classification and Business Customer Relationship Management in J. Phys. Conf. Ser. vol.1550.

[13]   A. José-García and W. Gómez-Flores, 2016 Automatic clustering using nature-inspired metaheuristic : A survey Appl. Soft Compt. vol.41 pp.192-213.

[14]   K. M. Passino, Biomimicry of bacterial foraging for distributed optimization and control, IEEE Control Systems Magazine (Volume: 22, Issue: 3, June 2002).

[15]   R. A. Ofosu, S.I. Kamau, J.N. Nderu, et al., Determination of Optimal PI Gains For Fuzzy-PI Controller Using Bacterial Foraging Algorithm (BFA), IOSR Journal of Electrical and 588 Electronics Engineering 11(2) (2016), 26–33.

[16]   D. Guo and J. Zhou, Numerical integration based on bacterial foraging algorithm, Science and Technology Vision 10 555 (2019), 118–120. 556.

[17]   Tripathy, M., Mishra, S., Lai, L.L., Zhang, Q.P.: Transmission loss reduction based on FACTS and bacteria foraging algorithm. In: Proceedings of PPSN, pp. 222–231 (2006).

[18]   Li, M.S., Tang, W.J., Tang, W.H., Wu, Q.H., Saunders, J.R.: Bacteria foraging algorithm with varying population for optimal power flow. In: Proceedings of EvoWorkshops 2007. LNCS, vol. 4448, pp. 32–41 (2007).

[19]   Biswas, A., Dasgupta, S., Das, S., Abraham, A.: Synergy of PSO and bacterial foraging optimization: a comparative study on numerical benchmarks. In: Proceedings 2nd International Symposium Hybrid Artificial Intelligent Systems (HAIS). Advances Soft Computing Series, Innovations in Hybrid Intelligent Systems. ASC, vol. 44, pp. 255–263. Springer, Germany (2007).

[20]   Korani, W.: Bacterial foraging oriented by particle swarm optimization strategy for PID tuning. In: GECCO'08 Proceedings of the Genetic and Evolutionary Computation Conference. ACM, pp. 1823–1826. Atlanta (2008).

[21]   Dasgupta, S., Biswas, A., Das, S., Panigrahi, B.K., Abraham, A.: A Micro-Bacterial Foraging Algorithm for High-Dimensional Optimization (2009).

[22]   Chen, H., Zhu, Y., Hu, K.: Cooperative bacterial foraging optimization. Discret. Dyn. Nat. Soc. 2009.

[23]   Dasgupta, S., Das, S., Abraham, A., Biswas, A.: Adaptive computational chemotaxis in bacterial foraging optimization: an analysis. IEEE Trans. Evolut. Comput. 13(4), 919–419(2009).

[24]   Chen, H., Zhu, Y., Hu, K.: Multi-colony bacteria foraging optimization with cell-to-cell communication for RFID network planning. Appl. Soft Comput. 10, 539–47 (2010).

[25]   Kim, D.H.: Hybrid GA-BF based intelligent PID controller tuning for AVR system. Appl. Soft Comput. 11, 11–22 (2011).

[26]   Gollapudi, S.V.R.S., Pattnaika, S.S., Bajpaib, O.P., Devi, S., Bakwad, K.M.: Velocity modulated bacterial foraging optimization technique (VMBFO). Appl. Soft Comput. 11, 154–65 (2011).

[27]   Okaeme, N.A., Zanchetta, P.: Hybrid bacterial foraging optimization strategy for automated experimental control design in electrical drives. IEEE Trans. Ind. Inf. 9, 668–8 (2013).

[28]   Abd-Elazim, S.M., Ali, E.S.: A hybrid particle swarm optimization and bacterial foraging for optimal power system stabilizers design. Electr. Power Energy Syst. 46, 334–41 (2013).

[29]   Mandeep Kaur ,Sanjay Kadam: A novel multi-objective bacteria foraging optimization algorithm (MOBFOA) for multi-objective scheduling, Applied Soft Computing , Volume 66, May 2018, Pages 183-195.

[30]   Lv, X.; Chen, H.; Zhang, Q.; Li, X.; Huang, H.; Wang, G. An Improved Bacterial-Foraging Optimization-Based Machine Learning Framework for Predicting the Severity of Somatization Disorder. Algorithms 2018, 11, 17.

[31]   Huang Chen, Lide Wang,Jun Di, and Shen Ping,: Bacterial Foraging Optimization Based on Self-Adaptive Chemotaxis Strategy, Computational Intelligence and Neuroscience, Volume 2020 | Article ID 2630104.

[32]   Yufang Dan, Jianwen Tao; Knowledge worker scheduling optimization model based on bacterial foraging algorithm, Future Generation Computer Systems, Volume 124,2021,Pages 330-337,ISSN 0167-739X.

[33]   Bo Yang, Xuelin Huang, Weizheng Cheng, Tao Huang, Xu Li, Discrete bacterial foraging optimization for community detection in networks, Future Generation Computer Systems, Volume 128,2022,Pages 192-204, ISSN 0167-739X.

[34]   Sandeep Gogula, V. S. Vakula, Optimization for position and rating of distributed generating units using bacteria foraging algorithm to reduce power losses, International Journal of Cognitive Computing in Engineering ,Volume 4,2023,Pages 287-300,ISSN 2666-3074,

[35]   C. Blake, E. Keough and C. J. Merz, UCI repository of machine learning database (1998). http://www.ics.uci.edu/~mlearn/MLrepository.html.

# Traffic Flow Prediction in Urban Networks: Integrating Sequential Neural Network Architectures

Eva Lieskovska, Maros Jakubec, Pavol Kudela

University Science Park, University of Zilina, Zilina, Slovak Republic

*Abstract*—The rapid growth of urban areas has significantly compounded traffic challenges, amplifying concerns about congestion and the need for efficient traffic management. Accurate short-term traffic flow prediction remains important for strategic infrastructure planning within these expanding urban networks. This study explores a Transformer-based model designed for traffic flow prediction, conducting a comprehensive comparison with established models such as Long Short-Term Memory (LSTM), Bidirectional Long Short-Term Memory (BiLSTM), Bidirectional Gated Recurrent Unit (BiGRU), and Time-Delay Neural Network (TDNN). Our approach integrates traditional time series values with derived time-related features, enhancing the model's predictive capabilities. The aim is to effectively capture temporal dependencies within operational data. Despite the effectiveness of existing models, internal complexities persist due to diverse road conditions that influence traffic dynamics. The proposed Transformer model consistently demonstrates competitive performance and offers adaptability when learning from longer time spans. However, the simpler BiLSTM model proved to be the most effective when applied to the utilized data.

*Keywords*—*Traffic flow; short-term prediction; machine learning; transformer*

## I. INTRODUCTION

Urbanization and the subsequent surge in vehicular traffic pose challenges to the efficiency and sustainability of urban transportation networks. The intricate interplay of dynamic factors, including population growth, urban expansion, and evolving commuter behaviours, necessitates innovative solutions for managing traffic flow. In particular, the advent of advanced predictive models has emerged as a cornerstone in addressing the complexities inherent in urban traffic dynamics [1], [2].

Traffic flow prediction, an important component of intelligent transportation systems, facilitates proactive traffic management, congestion alleviation, and resource optimization. This predictive capability is increasingly crucial in urban planning and policymaking. Precise insights into future traffic patterns empower decision-makers to devise effective strategies for infrastructure development, traffic routing, and overall enhancement of urban mobility. Understanding population behaviours within transport models, especially in relation to mode choice for trips, forms a critical aspect that can influence the precision and application of predictive traffic models [3].

In the domain of traffic flow prediction, the quest for accurate, adaptive, and efficient models has intensified, given

the important role of predictive systems in optimizing urban transportation networks. Traditional models, including Long Short-Term Memory (LSTM), Gated Recurrent Unit (GRU), Bidirectional Recurrent models and Convolutional Neural Networks (CNNs) such as Time-Delay Neural Network (TDNN), have significantly contributed to unravelling temporal dependencies within traffic data [4].

The introduction of the Transformer architecture [5], initially developed for Natural Language Processing (NLP) tasks, has paved the way for sequence modelling in various domains. Known for its distinctive attention mechanisms, this architecture revolutionises sequential data processing by employing self-attention mechanisms. This allows for a deeper comprehension of intricate relationships within sequences. The ability to discern temporal correlations has highlighted its potential application in traffic flow prediction models [6]–[8].

In this study, we undertake the task of forecasting future vehicle counts based on historical observations, with a specific focus on univariate traffic flow forecasting. The objective is to harness the capabilities of a transformer, which excels at discerning intricate traffic dynamics. The aim is to analyse how well it can decode complex traffic patterns by capturing time-related nuances and dependencies within the traffic data. Furthermore, we compare the performance of the transformer with the established neural networks for sequence modelling, such as LSTM, Bidirectional Long Short-Term Memory (BiLSTM), Bidirectional Gated Recurrent Unit (BiGRU), and TDNN. The incorporation of temporal features and the evaluation of distinct past observation intervals might yield additional insights for our analysis.

The paper is organized as follows: Section II provides an overview of existing traffic flow prediction models, Section III offers a detailed description of the prediction models, Section IV explains the experimental setup and evaluation methodologies, and Section V presents an analysis and comparative assessment of results. In Section VI, the discussion of the results is presented, and Section VII concludes with remarks that outline implications for future research.

## II. RELATED WORKS

The evolution of time series prediction has been marked by advancements in data analysis, machine learning, and computational power. Initially, time series prediction relied on statistical (or parametric) methods such as Autoregressive (AR) and Moving Average (MA) models [9], [10]. These models are the building blocks of the Autoregressive Integrated Moving

Average (ARIMA) model [11], which remains a common approach for time series prediction to date. These methods assumed that future values depend linearly on past observations and aimed to capture the underlying trends and patterns.

Another frequently used method is employing Kalman filtering [12], [13]. Due to dynamic traffic conditions and the nonlinear nature of traffic flow, parametric methods may struggle to effectively capture traffic features. As a result, there has been a shift in focus towards non-parametric machine learning methods in the field of traffic flow forecasting [4]. Decision trees, k-nearest neighbour (k-NN) [14], Support Vector Machines, and Neural Networks (NNs) started making their way into the domain. However, challenges remained in handling the temporal dependencies inherent in time series data.

The resurgence of interest in neural networks, particularly Recurrent Neural Networks (RNNs), marked a significant milestone. RNNs, with their ability to capture sequential dependencies, demonstrated improved performance in time series prediction tasks. LSTMs are a type of RNN, which have shown the ability to extract complex correlations in non-linear traffic data and capture long-term dependencies. The study conducted in [15] compares the LSTM architecture with models such as random walk, support vector regression, wavelet neural network, and the stacked autoencoder, emphasizing its favourable outcomes in short-term traffic flow prediction. The hybrid LSTM proposed in [16] optimizes its structure and parameters to adapt to various traffic scenarios. Comparative analysis reveals that the hybrid LSTM model outperformed other typical models (Fuzzy C-Means, Kalman filter and LSTM) in terms of prediction accuracy. This improvement in accuracy was achieved with only a marginal increase in processing time compared to LSTM model.

The performance comparison between LSTM and its simplified counterpart GRU, indicates that GRU outperformed LSTM when the past observation sequences were small [17]. On the other hand, LSTM performed better with more complex datasets and required the use of extended sequences to predict future traffic volume. A comparative analysis with benchmark models proposed in [18], including ARIMA, LSTM, BiLSTM, and GRU, indicates the superior performance of the BiGRU model. The bidirectional model utilizes preceding and succeeding time sequences to extract additional traffic flow information. Notably, deep learning methods, including Bi-GRU, outperformed the traditional ARIMA model in prediction accuracy, particularly during peak periods. However, the BiGRU model exhibited a slight lag in traffic flow prediction.

Recent advances in time series forecasting using Transformer models are gaining traction in the field of traffic forecasting. Known for their prowess in cross-sequence tasks, these models have been refined to predict temporal data, fundamentally transforming conventional methodologies by optimizing computing processes and capturing extensive dependencies. Cai et al. [6] focused on addressing spatio-temporal dependencies in traffic forecasting. Their Traffic Transformer architecture, inspired by the Transformer framework and Graph CNNs, adeptly managed periodicity, and spatial dependencies. It showcased superior performance with real-world traffic datasets. Reza et al. [7] introduced a multi-head attention-based transformer model for traffic flow forecasting. The model demonstrates greater efficiency in capturing prolonged traffic flow patterns compared to recurrent-based models. However, to achieve optimal performance, the proposed transformer required substantial amounts of training data. Existing studies predominantly concentrate on short-term predictions, creating a gap in long-term traffic forecasting research. Tedjopurnomo et al. [8] stress the significance of extending prediction to 24 hours for better congestion planning. To overcome limitations in current recurrent structure-based models for long-term traffic prediction, they introduce a modified Transformer model named TrafFormer, incorporating time and day embedding. Experimental results highlight the superior performance of their proposed model compared to existing hybrid neural network models.

## III. METHODS

This section delineates the intricacies of sequential neural network architectures—LSTM, BiLSTM, BiGRU, TDNN, and Transformer—applied in the domain of traffic flow prediction models.

### A. Long Short-Term Memory and Bidirectional Long Short-Term Memory

LSTM, an extension of a Vanilla RNNs, presents a robust architecture aimed at resolving the limitations of conventional RNNs in capturing long-range dependencies. Addressing the vanishing gradient problem inherent in RNNs, LSTM units incorporate a memory cell ($c_t$) that persists and evolves over time steps.

At each time step, LSTM units navigate through three gates: the forget gate ($f_t$), input gate ($i_t$), and output gate ($o_t$). These gates modulate the flow of information, orchestrating the update and retention of information within the cell state. The LSTM architecture is shown in Fig. 1.



Fig. 1. Graphical visualization of the functioning of the LSTM unit.

Each gate within the LSTM unit serves a distinctive function:

- Forget gate ($f_t$): This gate regulates the relevance of past information, allowing the LSTM unit to decide the degree of retention or discarding of prior information from the cell state.

- Input Gate ($i_t$): Responsible for modulating incoming information, the input gate enables the selective update of the cell state based on the present input sequence and the preceding state.

- Output gate ($o_t$): Governing the flow of information from the cell state to generate the output, this gate ensures the controlled dissemination of relevant information.

In the context of traffic flow prediction models, LSTM networks exhibit remarkable proficiency in capturing and predicting complex traffic dynamics over prolonged periods, owing to their capacity to capture long-term dependencies within sequential traffic data.

BiLSTM extends the capabilities of LSTM by incorporating bidirectional processing, allowing information to flow both forward and backward within the network. BiLSTM units consist of two LSTM layers: one processes the input sequence forward in time, while the other processes the sequence in reverse. Each BiLSTM unit operates with two sets of gates similar to LSTM: forget gates $(\overrightarrow{f_t}, \overleftarrow{f_t})$, input gates $(\overrightarrow{i_t}, \overleftarrow{i_t})$, and output gates $(\overrightarrow{o_t}, \overleftarrow{o_t})$ for the forward and backward directions, respectively. This dual directionality enables the network to capture dependencies in both past and future contexts simultaneously.

By leveraging information from both past and future contexts, BiLSTM units excel in comprehensively understanding the sequential nature of data. In the domain of traffic flow prediction models, BiLSTM architectures demonstrate enhanced capabilities in capturing complex temporal dependencies, leveraging bidirectional information flow to predict traffic patterns with improved accuracy, especially when dealing with nuanced traffic dynamics influenced by historical and future context [19].

### B. Gated Recurrent Unit and Bidirectional Gated Recurrent Unit

GRU presents an alternative architecture to LSTM, designed to capture long-range dependencies in sequential data. GRU units comprise two gates: reset gate ($r_t$) and an update gate ($z_t$), effectively regulating the flow of information within the network. The reset gate determines how much of the past information to forget, while the update gate modulates the blending of new input with the previous state. The GRU architecture is shown in Fig. 2.



Fig. 2. Graphical visualization of the functioning of the GRU unit.

Unlike LSTM, GRU units do not possess a separate cell state, simplifying the architecture while preserving its capacity to capture temporal dependencies. GRU units are adept at learning from sequential data due to their simplified structure, making them particularly suitable for traffic flow prediction models. Their ability to balance the preservation and update of past information allows for effective modelling of traffic dynamics, enabling the prediction of flow patterns with a focus on essential temporal relationships.

BiGRU extends the GRU architecture to process information bidirectionally. Similar to BiLSTM, BiGRU incorporates two sets of GRU layers that process input sequences in both forward and backward directions. BiGRU units maintain the characteristics of GRU but leverage bidirectional information flow, allowing simultaneous exploration of past and future contexts [20].

In this work, bidirectional RNNs were employed to improve training efficiency by simultaneously processing the input sequence in both forward and backward directions (see Table I). The BiGRU and BiLSTM models comprise two bidirectional recurrent layers, and their outputs are aggregated using global average pooling. For comparison, we also included the classical LSTM model, which comprises three sequential LSTM layers, an aggregating LSTM layer, and densely connected layers.

TABLE I. CONFIGURATION OF RECURRENT MODELS

| Layer | LSTM | BiLSTM | BiGRU |
|---|---|---|---|
| 1. | LSTM() | Bidirectional(LSTM) | Bidirectional(GRU) |
| 2. | Dropout(0.2) | Dropout(0.2) | Dropout(0.2) |
| 3. | LSTM() | Bidirectional(LSTM) | Bidirectional(GRU) |
| 4. | Dropout(0.2) | Dropout(0.2) | Dropout(0.2) |
| 5. | LSTM() | GlobalAvgPooling() | GlobalAvgPooling() |
| 6. | Dropout(0.2) | Dense() | Dense() |
| 7. | LSTM() | Dense(1) | Dense(1) |
| 8. | Dropout(0.2) | | |
| 9. | Dense() | | |
| 10. | Dense(1) | | |

### C. Time-Delay Neural Network

TDNN represents a specialized class of feedforward neural networks designed for modelling temporal sequences. These networks utilize fixed-size time windows to capture intricate temporal dependencies embedded within sequential data. Unlike recurrent counterparts such as LSTM or GRU, TDNNs employ distinct convolutional layers, each capturing unique temporal abstractions within the input data.

Operating through convolutional layers that traverse the input sequence, TDNNs adeptly extract features within predefined time windows or delays. These localized features then undergo further processing across subsequent layers, culminating in higher-level representations that encapsulate the temporal intricacies within the data. By focusing on local patterns across diverse time scales, TDNNs excel in capturing short and medium-term dependencies inherent in sequential data. The complete architecture of the TDNN used in our experiments is detailed in Table II.

TABLE II.        CONFIGURATION OF TDNN MODEL

| Layer | TDNN |
|-------|------|
| 1. | TDNNLayer([-2,2]) |
| 2. | TDNNLayer([-2,0,2]) |
| 3. | TDNNLayer([-3,0,3]) |
| 4. | TDNNLayer([0]) |
| 5. | TDNNLayer([0]) |
| 6. | Flatten() |
| 7. | Dense(32) |
| 8. | Dense(1) |

### D. Transformer

Transformers have emerged as a paradigm-shifting architecture within neural networks, initially recognized for their success in NLP tasks. Unlike traditional RNNs, Transformers process input data in parallel, disassembling it into smaller tokens embedded within high-dimensional vectors. These vectors are then passed through multiple layers, utilizing a mechanism called self-attention to focus on important input segments. This intrinsic mechanism empowers Transformers to capture long-range dependencies and effectively model the underlying structures of natural language. The utilization of Transformers in traffic flow prediction represents a frontier where their prowess in capturing contextual relationships and long-range dependencies can significantly contribute to the evolution of precise traffic flow prediction models.



Fig. 3.   Transformer architecture.

The Fig. 3 shows the proposed architecture of a single-block Transformer (where *Nx* represents the block ID), consisting of following sub-layers: a multi-head self-attention mechanism, an LSTM layer, and fully connected feed-forward network. In our implementation, we opted for the use of two transformer blocks based on experimental findings. The output of the last Transformer block is aggregated using row-wise and column-wise attention pooling and is then fed to the final dense layers. The model takes as input either the one-dimensional time series or two-dimensional time series $\times$ number of features.

## IV.   EXPERIMENTS

This section provides an overview of the experimental setup, dataset specifics, training strategies, and evaluation metrics crucial for both the development and assessment of the performance of the neural network architectures used in traffic flow prediction.

### A. Dataset

The traffic dataset [21] used in this study is publicly available on the Kaggle online platform. This dataset consists of a collection of time series data, recording vehicle counts at hourly intervals across four distinct junctions. The features within this dataset include DateTime, Junction Type, Vehicle Count, and ID. The temporal span of data collection varies, encompassing observations from November 2015 to June 2017 for three junctions and from January 2017 to June 2017 for the remaining junction. Overall, this dataset comprises a total of 48,100 observations, providing insights into the hourly vehicular traffic across multiple junctions. In this study, data from junction number one was selected for experimentation.

Preprocessing techniques, including Z-score normalization and differencing with a one-week window span, were employed to mitigate inherent temporal patterns and trends within the dataset. Normalization addresses the issue of diversity in the value ranges of time series data, which is suboptimal for neural network input. The stationarity of the data was assessed using the Augmented Dickey-Fuller test.

In addition to time series values, we also included derived time-related features such as the month, hour, day of the week, weekend indicator, and lag features representing the values from the previous hour and the same hour on the previous day.

### B. Network Setup and Training

The experiments were conducted on a hardware platform, encompassing the environmental parameters listed in Table III.

TABLE III.        EXPERIMENTAL SETUP

| Parameters | Configuration |
|------------|---------------|
| CPU | Intel Core i9-12900HX |
| GPU | nVidia GeForce RTX 3080 Ti |
| GPU memory size | 16GB |
| RAM | 64GB |
| Operating systems | Win11 |
| Deep learning architecture | Tensorflow 2.10.1 |

Training of the neural networks—LSTM, BiLSTM, BiGRU, TDNN, and Transformer—entailed parameter tuning. These models were systematically constructed with iterative exploration into diverse epochs, learning rates, batch sizes, and optimizer choices. Furthermore, the Halving Grid Search algorithm was used to narrow down the search for optimal settings through successive halving. The key parameters governing model training are detailed in Table IV.

TABLE IV.        KEY PARAMETERS DURING MODEL TRAINING

| Parameters | Setup |
|---|---|
| Epochs | 500 |
| Early stopping patience | 10 |
| Momentum | 0.99 |
| Learning rate | 0.001 |
| Weight decay | 0.0005 |
| Batch size | 128 |
| Optimizer | Adam/Lion |

*C. Metrics*

The evaluation metrics are important in assessing the efficacy of traffic flow prediction models developed using neural networks. While analytical or theoretical validation of these models proves challenging, error metrics play a crucial role in assessing their performance [22].

The evaluation metric used in this work is the Mean Squared Error (MSE) and Mean Absolute Error (MAE). MSE is a common metric employed to measure the average squared difference between the actual and predicted values (1). A higher MSE indicates greater prediction error.

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_{ref_i} - y_{pred_i})^2 \qquad (1)$$

where, n denotes the number of values. In this study, the Root Mean Square Error (RMSE) was utilized, which is the square root of MSE. This choice was made because RMSE shares the same scale as the original target variable.

The squaring of deviations in MSE significantly impacts the results, especially for extreme values. MSE exhibits higher sensitivity to these outliers. Conversely, for proximal values, squaring produces even smaller values, indicating their reduced significance rendering MSE less sensitive to nearby values. Therefore, an additional metric was employed to assess the performance of the models.

MAE operates similarly to MSE and represents the average positive deviation between predicted values and reference values. MAE is computed as the average absolute difference between predicted and reference values in Eq. (2) for n instances:

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_{ref_i} - y_{pred_i}| \qquad (2)$$

MAE provides a single value encapsulating all absolute deviations. The MAE metric does not amplify the effect of outliers since it considers absolute differences without squaring.

## V.    EXPERIMENTAL RESULTS

In this section, we present and analyse the experimental results obtained from applying various time-series forecasting models. The outcomes of the model evaluation are detailed in Table V, encompassing results for five distinct models: LSTM, BiLSTM, BiGRU, TDNN, and Transformer. For each model, performance is assessed across two time intervals—6 hours and 12 hours. Past observations from the last t hours served as input, and predictions for the subsequent time point (t + 1 hour) were generated. Furthermore, two experimental settings were employed: time series modelling using simple sequences, and time series modelling with additional features. During the evaluation phase, we conducted 10 successive model trainings, and the results of the best model are reported.

The results reveal variations in the models' predictive capabilities under different forecasting horizons. In most cases, models that make predictions based on the past 6-hour time interval achieved better results. The Transformer model appears to perform well when forecasting based on longer time spans. The complexity of the proposed Transformer might handle intricate inputs more efficiently.

In general, the inclusion of time features in the learning process resulted in improved error metrics. By integrating temporal information, models acquire the capability to leverage inherent temporal patterns and dependencies within time series data. The best results for each time interval (columns) are highlighted in bold. Among the considered models, BiLSTM, BiGRU, and the proposed Transformer proved to be the most effective, with BiLSTM achieving the highest performance. This outcome can be attributed to the fact that a simpler model is more suitable for a smaller database. While the proposed Transformer model consistently demonstrates competitive performance, particularly evident with MAE values ranging from 0.1695 to 0.1714, it may be better suited for larger datasets. The utilization of pretraining could potentially further enhance its performance.

TABLE V.        COMPARISON OF THE EXPERIMENTAL RESULTS

| Model | Metrics | Time series | | Time series × features | |
|---|---|---|---|---|---|
| | | 6h | 12h | 6h | 12h |
| LSTM | RMSE | 0.2388 | 0.2399 | 0.2365 | 0.2383 |
| | MAE | 0.1720 | 0.1725 | 0.1694 | 0.1703 |
| BiLSTM | RMSE | 0.2380 | 0.2392 | 0.2350 | 0.2363 |
| | MAE | 0.1708 | 0.1717 | 0.1688 | 0.1692 |
| BiGRU | RMSE | 0.2361 | 0.2398 | 0.2359 | 0.2368 |
| | MAE | 0.1704 | 0.1715 | 0.1692 | 0.1699 |
| TDNN | RMSE | 0.2386 | 0.2406 | 0.2388 | 0.2394 |
| | MAE | 0.1720 | 0.1735 | 0.1724 | 0.1727 |
| Transformer | RMSE | 0.2385 | 0.2367 | 0.2363 | 0.2376 |
| | MAE | 0.1714 | 0.1711 | 0.1711 | 0.1695 |

Fig. 4.   Five-day prediction comparison across various sequence models: a) simple time series prediction; b) time series prediction with additional features.

## VI.   Discussion

The evaluation of LSTM, BiLSTM, BiGRU, TDNN, and a modified Transformer over two time intervals (6 hours and 12 hours) and across two experimental settings, including time series modelling with simple sequences and time series modelling with additional features, has provided insights into their predictive capabilities. The effectiveness of models is influenced by the choice of the past time horizon. Notably, most of the models learning from a 6-hour time span demonstrated superior performance compared to those learning from 12 hours. The inherent complexity of the Transformer architecture enables it to effectively capture temporal dependencies, making it particularly well-suited for forecasting based on longer time spans.

The predicted outcomes of all models without the use of time features are visualized in Fig. 4(a). Upon comparison with the addition of time features in Fig. 4(b), subtle improvements in prediction accuracy can be observed. The visualized days start from Tuesday and extend until midday on Sunday. The predicted values closely mimic the real-world values, with one notable exception: the models have learned to anticipate an increase in the number of vehicles on Thursdays and Fridays. Including supplementary information about holidays or non-working days might improve the model's decision-making process, especially in pinpointing the busiest traffic days of the week related to holiday travel.

The prediction outcomes for the proposed Transformer model and the best performing BiLSTM are illustrated in Fig. 5. For a more detailed perspective, only two days are displayed, revealing a distinct decline in the number of vehicles from Friday to Saturday. The incorporation of temporal features (see Fig. 5(b)) to some extent helped align the predicted values more closely with the actual values.

Selecting between BiLSTM and Transformer for time series prediction relies on the characteristics of the data and the available computational resources. While BiLSTM is a type of

RNN that can capture temporal dependencies in sequential data, Transformer is a type of attention-based NN that can process sequential data in parallel, resulting in faster training times. The size of a dataset can influence the performance difference between a BiLSTM and a Transformer. The Transformer can be well-suited for transfer learning, particularly when pre-trained on large datasets, making it valuable for tasks involving limited labelled data.



Fig. 5.   The prediction outcomes of BiLSTM and Transformer models: a) simple time series prediction; b) time series prediction with additional features.

## VII.   Conclusion

In this study, we conducted an analysis of various time-series forecasting models, including LSTM, BiLSTM, BiGRU, TDNN, and modified Transformer. The evaluation encompassed two time intervals (6 hours and 12 hours) and two experimental settings: time series modelling using simple sequences and time series modelling with additional features. Our findings indicate variations in the predictive capabilities of the models under different forecasting horizons. Notably, models learning from a 6-hour time interval generally outperformed those learning from 12 hours. The Transformer model demonstrated efficacy in longer time spans, showcasing its ability to handle intricate inputs efficiently due to its inherent complexity.

The integration of time features into the learning process often resulted in improvements in error metrics. This enhancement arises from the models' capacity to leverage temporal patterns within time series data. Among the considered models, BiLSTM, BiGRU, and the proposed Transformer emerged as the most effective, with BiLSTM achieving the highest performance.

In our future work, the potential of transfer learning and improved fine-tuning will be explored. Moreover, evaluating other time series datasets may provide additional insights into the proposed analysis. The findings of this study can contribute to the broader understanding of model selection and optimization in time series forecasting, with implications for both research and practical applications in urban planning and traffic management systems.

## REFERENCES

[1] A. Boukerche and J. Wang, 'Machine Learning-based traffic prediction models for Intelligent Transportation Systems', Comput. Netw., vol. 181, p. 107530, Nov. 2020, doi: 10.1016/j.comnet.2020.107530.

[2] R. S. Joshi et al., 'State-of-the-art reviews predictive modeling in adult spinal deformity: applications of advanced analytics', Spine Deform., vol. 9, no. 5, pp. 1223–1239, Sep. 2021, doi: 10.1007/s43390-021-00360-0.

[3] M. Cingel, M. Drliciak, J. Celko, K. Zabovska, 'Modal Split Analysis by Best-Worst Method And Multinominal Logit Model.', presented at the Transport Problems: an International Scientific Journal . 2023, Vol. 18 Issue 1, p55-65. 11p.

[4] D. A. Tedjopurnomo, Z. Bao, B. Zheng, F. M. Choudhury, and A. K. Qin, 'A Survey on Modern Deep Neural Network for Traffic Prediction: Trends, Methods and Challenges', IEEE Trans. Knowl. Data Eng., vol. 34, no. 4, pp. 1544–1561, Apr. 2022, doi: 10.1109/TKDE.2020.3001195.

[5] K. He, X. Zhang, S. Ren, and J. Sun, 'Spatial Pyramid Pooling in Deep Convolutional Networks for Visual Recognition', in Computer Vision – ECCV 2014, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds., in Lecture Notes in Computer Science. Cham: Springer International Publishing, 2014, pp. 346–361. doi: 10.1007/978-3-319-10578-9_23.

[6] L. Cai, K. Janowicz, G. Mai, B. Yan, and R. Zhu, 'Traffic transformer: Capturing the continuity and periodicity of time series for traffic forecasting', Trans. GIS, vol. 24, no. 3, pp. 736–755, 2020, doi: 10.1111/tgis.12644.

[7] S. Reza, M. C. Ferreira, J. J. M. Machado, and J. M. R. S. Tavares, 'A multi-head attention-based transformer model for traffic flow forecasting with a comparative analysis to recurrent neural networks', Expert Syst. Appl., vol. 202, p. 117275, Sep. 2022, doi: 10.1016/j.eswa.2022.117275.

[8] D. A. Tedjopurnomo, F. M. Choudhury, and A. K. Qin, 'TrafFormer: A Transformer Model for Predicting Long-term Traffic'. arXiv, Mar. 02, 2023. doi: 10.48550/arXiv.2302.12388.

[9] S. Xu and B. Zeng, 'Network Traffic Prediction Model Based on Auto-regressive Moving Average', J. Netw., vol. 9, no. 3, pp. 653–659, Mar. 2014, doi: 10.4304/jnw.9.3.653-659.

[10] M.-C. Tan, S. C. Wong, J.-M. Xu, Z.-R. Guan, and P. Zhang, 'An Aggregation Approach to Short-Term Traffic Flow Prediction', IEEE Trans. Intell. Transp. Syst., vol. 10, no. 1, pp. 60–69, Mar. 2009, doi: 10.1109/TITS.2008.2011693.

[11] N. L. Nihan and K. O. Holmesland, 'Use of the box and Jenkins time series technique in traffic forecasting', Transportation, vol. 9, no. 2, pp. 125–143, Jun. 1980, doi: 10.1007/BF00167127.

[12] J. Guo, W. Huang, and B. M. Williams, 'Adaptive Kalman filter approach for stochastic short-term traffic flow rate prediction and uncertainty quantification', Transp. Res. Part C Emerg. Technol., vol. 43, pp. 50–64, Jun. 2014, doi: 10.1016/j.trc.2014.02.006.

[13] S. V. Kumar, 'Traffic Flow Prediction using Kalman Filtering Technique', Procedia Eng., vol. 187, pp. 582–587, Jan. 2017, doi: 10.1016/j.proeng.2017.04.417.

[14] Y. Chen, Y. Zhang, and J. Hu, 'Multi-Dimensional traffic flow time series analysis with self-organizing maps', Tsinghua Sci. Technol., vol. 13, no. 2, pp. 220–228, 2008, doi: 10.1016/S1007-0214(08)70036-1.

[15] H. Shao and B.-H. Soong, 'Traffic flow prediction with Long Short-Term Memory Networks (LSTMs)', in 2016 IEEE Region 10 Conference (TENCON), Nov. 2016, pp. 2986–2989. doi: 10.1109/TENCON.2016.7848593.

[16] Y. Xiao and Y. Yin, 'Hybrid LSTM Neural Network for Short-Term Traffic Flow Prediction', Information, vol. 10, no. 3, Art. no. 3, Mar. 2019, doi: 10.3390/info10030105.

[17] L. C. Das, 'Traffic Volume Prediction using Memory-Based Recurrent Neural Networks: A comparative analysis of LSTM and GRU'. arXiv, Mar. 22, 2023. doi: 10.48550/arXiv.2303.12643.

[18] S. Wang, C. Shao, J. Zhang, Y. Zheng, and M. Meng, 'Traffic flow prediction using bi-directional gated recurrent unit method', Urban Inform., vol. 1, no. 1, p. 16, Dec. 2022, doi: 10.1007/s44212-022-00015-z.

[19] R. L. Abduljabbar, H. Dia, and P.-W. Tsai, 'Unidirectional and Bidirectional LSTM Models for Short-Term Traffic Prediction', J. Adv. Transp., vol. 2021, p. e5589075, Mar. 2021, doi: 10.1155/2021/5589075.

[20] C. Chai et al., 'A Multifeature Fusion Short-Term Traffic Flow Prediction Model Based on Deep Learnings', J. Adv. Transp., vol. 2022, p. e1702766, May 2022, doi: 10.1155/2022/1702766.

[21] Fedesoriano, Traffic Prediction Dataset, February 2021.Retrieved from https://www.kaggle.com/datasets/fedesoriano/traffic-prediction-dataset.

[22] N. A. M. Razali, N. Shamsaimon, K. K. Ishak, S. Ramli, M. F. M. Amran, and S. Sukardi, 'Gap, techniques and evaluation: traffic flow prediction using machine learning and deep learning', J. Big Data, vol. 8, no. 1, p. 152, Dec. 2021, doi: 10.1186/s40537-021-00542-7.

# Experience Replay Optimization via ESMM for Stable Deep Reinforcement Learning

Richard Sakyi Osei, Daphne Lopez

School of Computer Science Engineering and Information Systems
Vellore Institute of Technology, Vellore, India

*Abstract*—The memorization and reuse of experience, popularly known as experience replay (ER), has improved the performance of off-policy deep reinforcement learning (DRL) algorithms such as deep Q-networks (DQN) and deep deterministic policy gradients (DDPG). Despite its success, ER faces the challenges of noisy transitions, large memory sizes, and unstable returns. Researchers have introduced replay mechanisms focusing on experience selection strategies to address these issues. However, the choice of experience retention strategy has a significant influence on the selection strategy. Experience Replay Optimization (ERO) is a novel reinforcement learning algorithm that uses a deep replay policy for experience selection. However, ERO relies on the naïve first-in-first-out (FIFO) retention strategy, which seeks to manage replay memory by constantly retaining recent experiences irrespective of their relevance to the agent's learning. FIFO sequentially overwrites the oldest experience with a new one when the replay memory is full. To improve the retention strategy of ERO, we propose an experience replay optimization with enhanced sequential memory management (ERO-ESMM). ERO-ESMM uses an improved sequential retention strategy to manage the replay memory efficiently and stabilize the performance of the DRL agent. The efficacy of the ESMM strategy is evaluated together with five additional retention strategies across four distinct OpenAI environments. The experimental results indicate that ESMM performs better than the other five fundamental retention strategies.

*Keywords—Experience replay; experience replay optimization; experience retention strategy; experience selection strategy; replay memory management*

## I. INTRODUCTION

Deep reinforcement learning (DRL) has emerged as a robust framework for training agents to make intelligent decisions in complex environments [1] in the area of health [2], [3], [4], automobile [5], [6], robotics [7], [8], energy [9], [10], [11] and others. RL algorithms aim to maximize cumulative rewards by allowing an agent to interact with an environment, learning from trial and error. One crucial component of RL is experience replay (ER), which involves reusing past experiences to enhance the learning process. ER has been successful in the RL field since its implementation in the deep Q-network (DQN) algorithm [12] and has undergone numerous improvements by researchers. It has proven to be a valuable technique, facilitating improved sample efficiency, breaking transition correlations, stabilizing learning dynamics, and reducing the cost of training [12], [13], [14], [15]. However, the effectiveness of ER strongly relies on selecting and retaining relevant experiences [16]. Experience sampling

involves the sequential or stochastic (random) selection of an experience index (the array index of the stored transitions) to determine the experience to use for training the RL agent. In contrast, experience retention focuses on strategies that can be adopted to determine the experiences to be stored and how these experiences can be managed in the replay buffer to optimize the learning process.

In recent years, various ER strategies have been proposed to address the challenges associated with experience retention. Isele and Cosgan [17] suggests that strategies for sampling experiences can be based on surprise, reward, state-space coverage, global training distribution, and state-action similarities. Similarly, de Bruin et al. [16] identified the full database (Full DB), first-in-first-out (FIFO), temporal difference error (TDE), and exploration as strategies for experience retention. They assert that while the FIFO strategy sequentially replaces old experiences irrespective of their relevance to learning, the exploration strategy stochastically overwrites the least-explored experience. The TDE strategy stochastically overwrites the least surprising experience, while the Resv approach ensures that "observed" experiences have equal retention rights. Moreover, the Full DB strategy uses a large amount of memory to store all experiences and does not require the removal of experiences from the replay memory.

Even with the improvements in ER strategies, they have limitations. The TDE is susceptible to noise, differences in function approximation accuracy, and randomness in the environment [15], [16], [18], [19], [20]. The FIFO strategy is affected by rapid changes in the state distribution, and Resv has premature data distribution convergence and poor coverage of state action over the optimal policy. Exploration takes longer to learn the correct value function, and Full DB promotes the rise of irrelevant experiences in the replay buffer [16]. Hence, there is still ample room for improvement in designing more efficient and effective experience retention strategies.

The proposed experience replay optimization with enhanced sequential memory management (ERO-ESMM) is a novel reinforcement learning algorithm that aims to improve the learning stability of DRL agents. The ERO-ESMM algorithm uses an enhanced sequential memory management (ESMM) strategy to manage the replay memory efficiently and stabilize the agent's performance. Compared to five existing experience retention strategies, the experimental results indicate that ERO-ESMM exhibits superior performance.

In this study, we present an enhanced experience retention strategy for DRL. Our proposed strategy aims to improve the

efficiency and effectiveness of experience replay by carefully managing experiences within the replay memory. The primary advancements presented in this study can be outlined as follows:

*1)* Firstly, we develop three new retention strategies to improve the efficiency and effectiveness of experience replay.

*2)* Secondly, we investigate the effects of six retention strategies, including the enhanced FIFO, on the ERO-enhanced DDPG algorithm.

*3)* Finally, we propose an enhanced framework incorporating the highest-performing retention strategy into the ERO framework.

We review existing experience retention strategies in Section III to achieve these objectives. We then present our enhanced retention strategy, outlining its core principles and rationale in Section III. Subsequently, we describe the experimental setup used to evaluate the performance of our approach in Section IV and provide detailed results and analysis. Finally, in Section V, we conclude our study and recommend future work in the field.

## II. RELATED WORK

We recognize that sampling and retention strategies are essential to experience replay in reinforcement-learning algorithms. This section briefly overviews the experience replay mechanism and sampling strategy. However, the section focuses on retention strategies, explicitly identifying those that effectively improve the performance of RL algorithms that use experience replays. The section explains the actor-critic method since the study seeks to improve an actor-critic algorithm (DDPG).

### A. Experience Replay

With ER, an agent generates experiences using an exploration-exploitation method at specific intervals and stores them in fixed-size replay memory. The agent then samples these experiences uniformly and randomly from the replay buffer into a mini-batch and repeatedly uses them to train the RL algorithm. This random selection prevents high correlation among the sampled experiences.

ER was first introduced by Lin [21] in the early 1990s, but it gained widespread attention when Mnih et al. [12] combined it with a deep convolutional neural network to create a groundbreaking deep Q-network (DQN) algorithm. Since then, ER has become a critical component of RL algorithms, allowing them to use experience effectively and reduce the interactions required with the environment. Before the introduction of ER, algorithms such as Q-learning [20], which relied on a tabular data storage mechanism, could not retain previous state-action values because the current ones of the same state-action pair would overwrite them. This "catastrophic forgetting' [22], [23], [24], [25], [26] behavior leads to slower learning and poor algorithm convergence.

Although many algorithms that implement diverse sampling strategies have been developed, the size and data structure of the replay buffer, mini-batch size, experience retention rate, experience sampling, and retention techniques

significantly influence the performance of these algorithms [40]. The selection of the experience index for experience sampling or retention can be sequential or random (uniform or prioritized probability). Sequential index selection is not appropriate for experience sampling because it creates a high correlation among the selected experiences, which subsequently slows down the agent's learning [12], [13], [34], [41]. Table I shows that most cited RL algorithms use a sequential index selection approach to remove experiences from the replay memory.

TABLE I. SOME RL ALGORITHMS AND THEIR EXPERIENCE INDEX SELECTION APPROACH FOR EXPERIENCE RETENTION

| Algorithm | Experience Index Selection Approach | | |
|---|---|---|---|
| | *Sequential* | *Random* | |
| | | *Uniform* | *Priority* |
| Deep Q-Network (DQN) [12] | ✓ | | |
| Double DQN[27] | ✓ | | |
| Dueling DQN[28] | ✓ | | |
| Prioritized Experience Replay (PER) [29] | ✓ | | |
| Deep Deterministic Policy Gradient (DDPG) [30] | ✓ | | |
| Twin Delayed Deep Deterministic Policy Gradient (TD3) [18] | ✓ | | |
| Trust Region Policy Optimisation (TRPO) [31] | ✓ | | |
| Proximal Policy Optimisation (PPO) [32] | ✓ | | |
| Episodic Memory Deep Q-Network (EMDQN) [33] | ✓ | | |
| Advantage Actor-Critic (A2C) + Prioritized Stochastic Memory Management (PSMM) [20], [34] | | | ✓ |
| DQN + Dual Memory Structure (DMS) [12], [35] | ✓ | | ✓ |
| DDPG + Experience Replay Optimisation (ERO) [30], [36] | ✓ | | |
| Combined Experience Replay (CER) [37] | ✓ | | |
| Attentive Experience Replay (AER) [38] | ✓ | | |
| Selective Experience Replay (SER) [17] | | ✓ | ✓ |
| Prioritized Sequence Experience Replay (PSER) [39] | ✓ | | |

The advancements in ER are incorporated in many popular RL algorithms, such as DQN [12], dueling DQN [28], double DQN [27], twin delayed deep deterministic policy gradient (TD3)[18], deep deterministic policy gradient (DDPG) [30], proximal policy optimization (PPO)[32], episodic memory deep Q-Network (EMDQN) [33], and trust region policy optimization (TRPO) [31], still use the naïve ER uniform random sampling strategy. Other algorithms, such as prioritized experience replay (PER) [36], prioritized sequence experience replay (PSER) [37], experience replay optimization (ERO) [38], and attentive experience replay (AER) [39], implement prioritized strategies. Equally, prioritized stochastic memory management (PSMM) [20], combined experience replay (CER) [37], selective experience replay (SER) [17], and episodic memory control (EMC)[40] use experience retention

strategies (memory management strategies). In contrast, some replay strategies focus on the structure of the replay memory instead of the content[35], [42], [43]. ERO has proven superior among prioritized selection algorithms, owing to its easy adaptation and generalization to multiple environments [23].

### B. Experience Retention Strategies and Algorithms

Experience retention plays a critical role in the success of ER algorithms. We can only select the experiences available in the replay buffer for training. If valuable experiences are maintained in the buffer, there will be a higher probability of sampling a mini-batch full of relevant experience to train the RL agent. In contrast, the worst training could be given to the agent. Therefore, it is imperative to investigate and unearth innovative ways to improve existing retention strategies or, better still, develop new ones.

The naïve approach of randomly selecting buffered experiences uniformly or managing the replay memory with a simple FIFO strategy is simple but less successful than the prioritized approach for managing the replay buffer [16], [29], [44]. Recent enhanced works on experience replay have relied on rule-based strategies that directly prioritize transitions through sampling strategies or indirectly through retention strategies [36]. However, some prioritized strategies incorporate a certain degree of randomness during implementation, using hyperparameters to regulate prioritization. Prioritization relies on features such as the temporal difference error (TDE), reward signal, similarities or diversity of states[18], or a combination of any of these features [20], [36]. A comprehensive study by de Bruin et al. [16] outlines age, exploration, and surprise as the criteria for retaining experiences in the buffer.

Retention depends on the duration for which an experience remains in the buffer. FIFO, Full DB, and Reservoir are strategies that rely on age. Although FIFO uses sequential indexing to remove old experiences without regard for their contribution to learning, Reservoir overwrites experiences in a uniformly random fashion and has limited state-action coverage. The Full DB method accommodates all experiences until the end of the training but may retain irrelevant experiences.

It is worth noting that there are better choices than the exploration criteria when dealing with problems that require minimal interaction with the agent's environment [16]. TDE is an expression of the surprise between the targeted and predicted q-values. Overfitting can occur if not parameterized and regulated [29], [40]. Actor-Critic Method Two major approaches in RL, value-based and policy-based methods, have been widely explored. Value-based methods estimate the value function, while policy-based methods directly optimize the agent's policy [1]. However, each approach has its limitations. Value-based methods tend to suffer from overestimation, the curse of dimensionality, and are often computationally expensive. On the other hand, policy-based methods can be inefficient in exploring the environment and may need help with convergence [13], [18], [45]. To address these challenges, the actor-critic algorithm, a hybrid approach, combines the strengths of value-based and policy-based methods [13], [18], [46]. It consists of two key components: the actor and the critic. The actor represents the policy and selects actions based on the observed states. The critic estimates the value function, providing feedback to the actor by evaluating the chosen actions.

The actor is typically implemented as a parametric model, like a neural network, which maps states to a probability distribution over actions. It explores the environment, collects experiences, and adjusts its policy based on the rewards it receives. The critic, represented by a parametric model, estimates the value function by approximating the expected cumulative reward associated with states or actions. By combining the strengths of both approaches, actor-critic algorithms can achieve faster convergence, more stable learning, and better performance in a wide range of RL problems. This hybrid approach has found successful applications in various domains, including robotics control, game playing, and natural language processing [15], [34].

### III. METHODOLOGY

This section briefly introduces ERO and PSMM while paying particular attention to the selection and retention strategies used. It further presents the proposed framework and implemented algorithms.

### A. Experience Replay Optimization

ERO is an experience selection method that relies on Reward and TDE for prioritization [36]. Unlike other TDE prioritization sampling strategies that favor experiences with higher TD errors, ERO selects less surprised TDE experiences and uses a novel replay policy network for the prioritization process. A mini-batch of high-priority transitions (transitions with vector 1) was created, and its elements were uniformly sampled to train the agent [23]. After the agent interacts with the environment, the transitions are stored in the replay buffer and subsequently prioritized through a Boolean (0, 1) vectorization process using the replay policy. During training, the replay policy receives feedback from the environment for policy evaluation.

Since the performance of a sampling strategy is highly dependent on the implementing algorithm and the benchmark environment [16], there is the need for a sampling method that can learn and adapt to different algorithms and environments - ERO does rightly so. ERO still uses the FIFO retention strategy despite its novel adapting strategy and superior performance over prioritized sampling methods such as PER [36], [44].

Hence, when the replay memory exceeds its capacity, the oldest transition is sequentially replaced with a new transition, irrespective of its importance in learning. Nonetheless, when relevant transitions are retained and frequently sampled using an intelligent index selection strategy, we are optimistic that the agent's performance and convergence rate will improve [14], [45], [47]. Therefore, there is a need to augment ERO with a memory management mechanism that is better than FIFO[16]. The beauty and novelty of the ERO algorithm depend on its replay policy network, which relies on Eq. (1) to Eq. (4). Table II explains the notations used in the equations, and the replay policy update is presented in Algorithm 1 [36].

$$\lambda = \left\{ \phi\left(f_{B_i}|\theta^{\phi}\right)\middle| B_i \in B \right\} \in \mathbb{R}^N \tag{1}$$

$$B^s = \{B_i | B_i \in B \wedge I_i = 1\} \qquad (2)$$

$$r^r = r_\pi^c - r_{\pi'}^c \qquad (3)$$

$$P(I|\phi) = \sum_{j:B_j \in B^{batch}} r^r \nabla \theta^\phi [I_i \log \phi + (1 - I_i) \log(1 - \phi)] \qquad (4)$$

where, $\phi$ denotes the function approximator, $B_i$ is a transition in the replay buffer $B$. $\theta^\phi$ denotes the parameters of $\phi$, $f_{B_i}$ is a feature vector, and N is the number of transitions in a mini-batch. The priority score function is expressed as $\phi(f_{B_i}|\theta^\phi) \in (0,1)$, where the priority score is represented by Lambda $(\lambda)$. $I_i \in \{0,1\}$ is the Bernoulli distribution of sample $B^s$. The replay reward, cumulative reward of the current policy, and cumulative reward of the previous policy are denoted by $r^r$, $r_\pi^c$, and $r_{\pi'}^c$ respectively.

---

**Algorithm 1:** PolicyUpdate

---

Input:

    Cumulative reward of current policy $r_\pi^c$

    Cumulative reward of previous policy $r_{\pi'}^c$

Output:

    Sample subset $B^s$

Calculate replay reward based on (3)

For (each replay updating step) do

    | Randomly sample batch $\{B_i\}$ form B

    | Update replay policy based on (4)

End

Sample subset $B^s$ from B using (2)

---

### B. Prioritized Stochastic Memory Management

Experience replay selection strategies are often the focus of researchers. However, it is essential to note that a poorly designed retention technique can negatively impact the performance of the learning agent [16], [17], [20]. One proposed method for effective replay memory management is prioritized stochastic memory management (PSMM), which was introduced by Ko and Chang [35]. The PSMM employs a stochastic approach to remove the history with the least TDE or return when the replay memory is full, using the probability computed in Eq. (5).

$$P_i = \frac{exp\left(-\rho(\alpha_{actor}R_{norm}^{(i)} + \alpha_{critic}TDE_{norm}^{(i)})\right)}{\sum_{i=1}^{c} exp\left(-\rho(\alpha_{actor}R_{norm}^{(i)} + \alpha_{critic}TDE_{norm}^{(i)})\right)} \qquad (5)$$

The computation of the probability for elimination $(P_i)$ in Kwon and Chang's method involves the utilization of historical information through return and temporal difference error (TDE) metrics [25]. These metrics were normalized to restrict their values from 0 to 1, facilitating unbiased evaluations and promoting stable memory management. The method employs several hyperparameters, such as $\alpha_{actor}$, $\alpha_{critic}$, and $\rho$, which are optimized for improved performance. The $\alpha_{actor}$ and $\alpha_{critic}$ determine the relative weights assigned to the actor and critic components, respectively, whereas $\rho$ represents a probability tuning parameter. This approach ensures the effective utilization of historical data and facilitates the optimization of the method's parameters.

The computations for $R_{norm}^{(i)}$ and $TDE_{norm}^{(i)}$ are shown in Eq. (6) and Eq. (7), respectively.

$$R_{norm}^{(i)} = \begin{cases} \frac{R^i - R_{min}}{R_{max} - R_{min}}, & if\ R_{max} \neq R_{min} \\ 0 & otherwise \end{cases} \qquad (6)$$

$$TDE_{norm}^{(i)} = \begin{cases} \frac{TDE^i - TDE_{min}}{TDE_{max} - TDE_{min}}, & if\ TDE_{max} \neq TDE_{min} \\ 0 & otherwise \end{cases} \qquad (7)$$

$$P_i = \frac{\left(R_{norm}^{(i)}\right)^\alpha}{\sum_{i=1}^{c} \left(R_{norm}^{(i)}\right)^\alpha} \qquad (8)$$

$$P_i = \frac{exp\left(-\rho(\alpha_{actor}R_{norm}^{(i)})\right)}{\sum_{i=1}^{c} exp\left(-\rho(\alpha_{actor}R_{norm}^{(i)})\right)} \qquad (9)$$

TABLE II.      SYMBOLS AND NOTATIONS USED IN THIS SECTION

| Notation | Explanation |
|---|---|
| $\phi$ | Function approximator |
| $B$ | Replay buffer |
| $B_i$ | A transition in the replay buffer B |
| $\theta^\phi$ | Parameters of the function approximator |
| $f_{B_i}$ | Feature vector |
| $\lambda$ | Priority score |
| $r^r$ | Cumulative reward |
| $r_\pi^c$ | Cumulative reward of current policy |
| $r_{\pi'}^c$ | Cumulative reward of previous policy |
| $B^s$ | A specified batch size of sampled transitions |
| $\rho$ | Probability tuning parameter |
| $\alpha$ | Parameter for tuning the probability of elimination |

### C. Proposed Framework

Researchers have recently harnessed and combined various algorithms' strengths to create resilient, stable, and generalized hybrid algorithms. Our proposed framework amalgamates the ERO framework and enhances the FIFO retention strategy.



Fig. 1. A preliminary experiment was conducted to identify the optimum experience retention ratio. When the replay memory capacity is reached, the buffered experiences are overwritten using a ratio. A ratio of 2:8 means 20% of the old experiences are sequentially overwritten with new experiences, and 80% are retained. However, for a ratio of 8:2, only 20% of the buffered experiences are retained.

Fig. 2. Proposed Framework: experience replay optimization with enhanced sequential memory management (ERO-ESMM). Transitions from the environment are stored in the replay buffer. Mini-batches from the transitions are vectorized, prioritized by the replay policy, and sampled uniformly at random to train the agent. After training, the replay policy receives feedback for policy evaluation. When the memory is full, ESMM ensures that transitions in the first half of the replay memory are sequentially overwritten with new ones.

The ESMM, PSMM($\alpha$), and PSMM($\rho$) retention strategies were developed to create the new framework. While PSMM($\alpha$) and PSMM($\rho$) use Eq. (8) and Eq. (9), respectively, ESMM extends the FIFO retention strategy by sequentially overwriting older transitions in the first half of the memory when the memory is full. The one-half was arrived at after preliminary experiments were conducted using different ratios of the replay buffer for experience retention. Fig. 1, which shows the preliminary experiment results, confirms that an even distribution of old and new experiences in the replay buffer enhances the performance of the RL agent.

The ESMM, PSMM($\alpha$), and PSMM($\rho$) strategies and FIFO, Full DB, and Resv were further investigated to ascertain their effects on the ERO-enhanced DDPG algorithm. The strategy with the highest mean return, ESMM, was incorporated into the ERO framework to propose an improved framework, an experience replay optimization with enhanced sequential memory management (ERO-ESMM). The proposed framework and its algorithm are shown in Fig. 2 and Algorithm 2, respectively.

| **Algorithm 2:** ERO-ESMM Enhanced DDPG |
|---|
| Initialize policy $\pi$, replay policy $\phi$ and buffer B |
| For (each iteration) do |
|     For (each time-step t) do |
|       Select action $a_t$ according to $\pi$ and state $s_t$ |
|       Execute action $a_t$ and observe $s_{t+1}$ and $r_t$ |
|       If (B is full) then |
|         If (index i+1 == ½ len(B) ) then |
|           i=0 |
|         Else |
|           i =(i+1) mod len(B) |
|         End |
|         Store transition($s_t$, $a_t$, $r_t$, $s_{t+1}$) at $B_i$ |
|       End |
|       If (episode is complete) then |
|         Calculate the cumulative reward $r_\pi^c$ |
|         If ( $r_{\pi'}^c \neq null$ ) then |

$$B^s = \text{PolicyUpdate}(r_{\pi'}^c, r_\pi^c, B)$$

End
    Set $r_{\pi'}^c \leftarrow r_\pi^c$
  End
End
For (each training step) do
    Uniformly sample a batch $\{B_i\}$ from $B^s$
    Update the critic of $\pi$
    Update the actor of $\pi$
    Update the target networks
    Update the $\lambda$ for each transition in $\{B_i\}$
End
End

### D. Setup of RL Environment

To ascertain the efficiency of the proposed framework, we conducted a series of experiments in the Pendulum-v0, MountainCarContinuous-v0, LunarLandarContinuous-v2, and the BipedalWalker-v3 environments [48] of the OpenAI Gym as a platform for the evaluation and analysis of results. Screenshots of these environments are shown in Fig. 3.



Fig. 3. Screenshots of two Classic Control (top row) and two Box2D (down row) environments from the OpenAI Gym. Fig. 3(a) and (b) represent the Pendulum-v0 and MountainCarContinuous-v0 environments, respectively. Fig. 3(c) and (d) represents the LunarLandarContinuous-v2 and the BipedalWalker-v3 environments respectively.

Pendulum-v0 presents a classical inverted pendulum swing-up problem, which demands that the agent persistently swing up the pendulum from an initial arbitrary position until it attains an upright position while its 3-dimensional observation space comprises angle, acceleration, and angular velocity, its action space is continuous, ranging between -2.0 (anti-clockwise torque) and 2.0 (clockwise torque). The agent's rewards depend on its actions and the associated state.

The MountainCarContinuous-v0 environment is another benchmark classical control environment that requires the RL agent to apply actions to a car to reach the top of a hill as quickly as possible. It is an extension of the classic "MountainCar-v0" environment but with continuous action space, making it more suitable for problems requiring continuous control. It consists of a 2-dimensional observation space of the car's position and velocity and a continuous action space between -1.0 and 1.0. The RL agent is negatively rewarded each time an action is taken until the car reaches the

top of the hill. Hence, the agent applies continuous efficient actions to overcome the car's inertia and climb the hill.

LunarLanderContinuous-v2 is an extension of the original Box2D discrete action space LunarLander-v2 environment, where the agent controls a lunar lander attempting to land on a designated landing pad on the moon's surface. In contrast to the discrete version, this environment allows the agent to apply a range of continuous actions (a throttle that varies between 0 and 1 and a rotation angle between -1 and 1) instead of selecting from a fixed set of discrete actions. It has an 8-dimensional observation space that includes information about the lander's position, velocity, orientation, angular velocity, whether the legs are touching the ground, and whether the lander has successfully landed. The RL agent receives positive rewards for moving closer to the landing pad and a significant positive reward for landing safely. Negative rewards are given for using fuel, and a slight negative reward is given for each time step.

BipedalWalker-v3 is similarly a Box2D environment where the agent controls a bipedal robot with four legs and must learn to make it walk and navigate through a complex terrain while avoiding obstacles. The challenge lies in learning a coordinated sequence of actions to control the robot's joints and achieve stable and efficient locomotion. This environment is represented by a 24-dimensional observation space, which contains information about the position, velocity, angle, angular velocity, and state of the joints, feet, and lidar sensors. The agent receives positive rewards for progressing forward and avoiding obstacles. However, negative rewards are given for using excessive torque or falling, encouraging the agent to discover stable and effective locomotion strategies.

The suitability of these environments for the DDPG algorithm is attributable to the availability of continuous tasks[30]. Thus, the selection of these environments was motivated by their compatibility with the requirements of the DDPG algorithm and their standard use in previous studies for evaluating RL algorithms [49].

*E. Parameter Setting*

Because the proposed framework extends the ERO-enhanced DDPG algorithm, the experimental configurations were based on the OpenAI DDPG stable baseline[50], and the hyperparameters for the sampling and retention methods were in line with ERO and PSMM, respectively. However, the memory size and number of time steps were adjusted during implementation.

In Table III, six notations and their explanations are clearly shown to facilitate comprehension of our visualizations. Aside from the Full DB, which has the memory size and number of time steps set to $2 \times 106$, the other five retention strategies were evaluated with a memory size of $1 \times 106$ and time steps of $2 \times 106$. Similar to PSMM[25], we set $\rho = 2.0$, $\alpha_{actor} = 0.5$ and $\alpha_{critic} = 0.5$ when implementing PSMM($\rho$). In the PSMM($\alpha$), we rely on an α value 0.6[21]. Other parameters for the evaluation of experiments were done on an Intel(R) Xeon(R) CPU E3-1220 v6 @ 3.00GHz(4CPUs) with 32GB RAM and a Windows Server 2012R2 operating system.

TABLE III. EXPERIENCE RETENTION STRATEGIES EVALUATED IN THE STUDY

| Notation | Explanation |
|---|---|
| FIFO | Sequentially overwrite old experiences with new ones. |
| ESMM | Replace experiences in the first half of the buffer with new ones. |
| Full DB | Experiences are not overwritten. The replay memory stores all experiences. |
| Resv | Experiences are uniformly, at random, overwritten with new ones. |
| PSMM(α) | Experiences are stochastically overwritten based on (8) |
| PSMM(ρ) | Experiences are stochastically overwritten based on (9) |

## IV. RESULTS AND DISCUSSION

This section examines the comparative effectiveness of the retention strategies within each environment. The metric for quantifying the performance of the evaluated strategies is the mean return. The return, also known as cumulative reward or cumulative return, is a measure of the overall success of the RL agent in achieving its goals. It is the summation of rewards received by the agent at each time step from the start of an episode until its termination[1], [51], [52]. The return is typically used to evaluate and compare the performance of different algorithms and policies in a specific reinforcement learning task. The mean return can be derived as follows:

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T \qquad (10)$$

$$Mean\ return = \frac{G_t}{number\ of\ episodes} \qquad (11)$$

where, $G_t$ is the return, $R_{t+1}$ is the reward at time step *t*, and *T* is the final time-step.

The performance of each retention strategy was evaluated based on the progressive average returns computed using Eq. (11). The results are illustrated in Fig. 4 to Fig. 7.

In the Pendulum environment, as shown in Fig. 4, the ESMM and PSMM(α) strategies exhibit superior performance compared to FIFO, while Full DB and PSMM(ρ) strategies perform worse. The Resv strategy exhibits limited effectiveness as it faces challenges in learning optimal policies to stabilize the pendulum. This struggle is reflected in a mean return of -1061.29.

Regarding the MountainCarContinuous-v0 environment, as indicated in Fig. 5, the ESMM strategy outperforms all other models, including FIFO. The Full DB and Resv strategies show some improvement over FIFO, while PSMM(α) and PSMM(ρ) perform worse with rewards of -6.05 and -6.46, respectively.

In the LunarLanderContinuous-v2 environment, as shown in Fig. 6, the ESMM, PSMM(α), and PSMM(ρ) strategies demonstrate superior performance compared to FIFO. The Full DB model performs worse than FIFO and PSMM(ρ), while the Resv model exhibits the poorest performance.

Likewise, in the BipedalWalker-v3 environment, the ESMM and PSMM(α) models perform better than the FIFO strategy. However, the Full DB and Resv models perform worse than ESMM and PSMM(α), with PSMM(ρ) showing the poorest performance among all models in this environment.

Fig. 4. Performance comparison of the six retention strategies evaluated on the Pendulum-v0 environment.



Fig. 5. Performance comparison of the six retention strategies evaluated on the MountainCarContinuous-v0 environment.



Fig. 6. Performance comparison of the six retention strategies evaluated on the LunarLandarContinuous-v2 environment.



Fig. 7. Performance comparison of the six retention strategies evaluated on the BipedalWalker-v3 environment.

The results indicate significant model performance variation across different OpenAI Gym environments. The ESMM strategy generally exhibits better and more stable performance across the experimented environments, which can be attributed to its fair distribution of experiences in the replay buffer. While the PSMM($\alpha$) model performs better in the Pendulum environment, the PSMM($\rho$) strategy performs better than other models, except ESMM, in the Lunar Lander Continuous-v2 environment.

Conversely, the Full DB model tends to perform worse in most environments because all experiences are stored, including the worst experiences obtained during the early stages of the RL agent's training. The Resv strategy's performance shows inconsistency across the environments due to its entirely random strategy for identifying the experience to remove when the replay buffer is full.

## V. CONCLUSION

In this study, we developed three new retention strategies to improve the effectiveness of experience replay. These strategies include ESMM, PSMM($\alpha$), and PSMM($\rho$). The development of these strategies is significant, as they provide new alternatives for managing the memory of an RL agent in a reinforcement learning setting. These strategies addressed specific challenges encountered in experience replay, such as the trade-off between memory usage and retention of relevant experiences.

This study also investigated the effects of six different retention strategies on the ERO-enhanced DDPG algorithm. These strategies include FIFO, Full DB, Resv, PSMM ($\alpha$), PSMM($\rho$), and our proposed method (ESMM). The results of this investigation are significant because they provide insights into the comparative performance of different retention strategies and help identify the most effective strategy. This information can guide the design of more efficient reinforcement learning algorithms.

Finally, we propose an enhanced framework incorporating the highest-performing retention strategy into the ERO framework. This enhanced framework, called experience-replay optimization with enhanced sequential memory management (ERO-ESMM), significantly contributes to reinforcement learning. By integrating the best-performing retention strategy, the framework offers a more optimized approach to experience replay, leading to improved performance in reinforcement-learning tasks.

Overall, the experimental results suggest that developing new retention strategies, combined with their investigation and incorporation into existing frameworks, can significantly improve the performance of reinforcement learning algorithms. These results have implications for future research and demonstrate the importance of exploring new techniques for optimizing reinforcement learning. In the future, we will use a separate neural network to predict the index of the experience to delete from the replay buffer when it exceeds its limit.

## REFERENCES

[1] R. S. Sutton and A. G. Barto, Reinforcement Learning: an Introduction. MIT press, 2018.

[2] B. Varghese and S. Krishnakumar, 'Fast Fractal Coding of MRI Images using Deep Reinforcement Learning', IJACSA, vol. 12, no. 4, 2021, doi: 10.14569/IJACSA.2021.0120492.

[3] W. E. Fathy and A. S. Ghoneim, 'A Deep Learning Approach for Breast Cancer Mass Detection', International Journal of Advanced Computer Science and Applications (IJACSA), vol. 10, no. 1, Art. no. 1, 31 2019, doi: 10.14569/IJACSA.2019.0100123.

[4] S. Luo, 'Lung Cancer Classification using Reinforcement Learning-based Ensemble Learning', International Journal of Advanced Computer Science and Applications (IJACSA), vol. 14, no. 8, Art. no. 8, 54/30 2023, doi: 10.14569/IJACSA.2023.01408120.

[5] E. Jacinto, F. Martinez, and F. Martinez, 'Navigation of Autonomous Vehicles using Reinforcement Learning with Generalized Advantage Estimation', International Journal of Advanced Computer Science and Applications (IJACSA), vol. 14, no. 1, Art. no. 1, 31 2023, doi: 10.14569/IJACSA.2023.01401103.

[6] I. Rasheed, F. Hu, and L. Zhang, 'Deep reinforcement learning approach for autonomous vehicle systems for maintaining security and safety using LSTM-GAN', Vehicular Communications, vol. 26, p. 100266, 2020.

[7] Y. Han et al., 'Deep reinforcement learning for robot collision avoidance with self-state-attention and sensor fusion', IEEE Robotics and Automation Letters, vol. 7, no. 3, pp. 6886–6893, 2022.

[8] G. E. Setyawan, P. Hartono, and H. Sawada, 'Cooperative Multi-Robot Hierarchical Reinforcement Learning', International Journal of Advanced Computer Science and Applications (IJACSA), vol. 13, no. 9, Art. no. 9, Dec. 2022, doi: 10.14569/IJACSA.2022.0130904.

[9] D. An, F. Cui, and X. Kang, 'Optimal scheduling for charging and discharging electric vehicles based on deep reinforcement learning', Frontiers in Energy Research, vol. 11, 2023, Accessed: Dec. 13, 2023. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fenrg.2023.1273820.

[10] H. Zhao et al., 'Combination optimization method of grid sections based on deep reinforcement learning with accelerated convergence speed', Frontiers in Energy Research, vol. 11, 2023, Accessed: Dec. 13, 2023. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fenrg.2023.1269854.

[11] V. L. Narayanan, 'Reinforcement learning in wind energy - a review', International Journal of Green Energy, vol. 0, no. 0, pp. 1–24, 2023, doi: 10.1080/15435075.2023.2281329.

[12] V. Mnih et al., 'Human-level control through deep reinforcement learning', Nature, vol. 518, no. 7540, pp. 529–533, 2015, doi: 10.1038/nature14236.

[13] Z. Wang et al., 'Sample efficient actor-critic with experience replay', in 5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings, 2017.

[14] D. Yang, X. Qin, X. Xu, C. Li, and G. Wei, 'Sample Efficient Reinforcement Learning Method via High Efficient Episodic Memory', IEEE Access, vol. 8, pp. 129274–129284, 2020, doi: 10.1109/ACCESS.2020.3009329.

[15] T. de Bruin, Sample efficient deep reinforcement learning for control. 2019. doi: 10.4233/uuid.

[16] T. De Bruin, J. Kober, K. Tuyls, and R. Babuška, 'Experience selection in deep reinforcement learning for control', Journal of Machine Learning Research, vol. 19, pp. 1–56, 2018.

[17] D. Isele and A. Cosgun, 'Selective experience replay for lifelong learning', in Proceedings of the AAAI Conference on Artificial Intelligence, 2018, pp. 3302–3309.

[18] S. Fujimoto, H. V. Hoof, and D. Meger, 'Addressing Function Approximation Error in Actor-Critic Methods', arXiv preprint arXiv:1802.09477, 2018.

[19] C. Kang, C. Rong, W. Ren, F. Huo, and P. Liu, 'Deep Deterministic Policy Gradient Based on Double Network Prioritized Experience Replay', IEEE Access, vol. 9, 2021, doi: 10.1109/ACCESS.2021.3074535.

[20] T. Kwon and D. E. Chang, 'Prioritized Stochastic Memory Management for Enhanced Reinforcement Learning', in 2018 IEEE International Conference on Consumer Electronics - Asia, ICCE-Asia 2018, 2018, pp. 7–10. doi: 10.1109/ICCE-ASIA.2018.8552124.

[21] L. Lin, 'Self-improvement Based On Reinforcement Learning, Planning and Teaching', in Machine Learning Proceedings 1991, vol. 321, 1992, pp. 323–327. doi: 10.1016/b978-1-55860-200-7.50067-2.

[22] Z. Li and D. Hoiem, 'Learning without forgetting', Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), vol. 9908 LNCS, pp. 614–629, 2016, doi: 10.1007/978-3-319-46493-0_37.

[23] J. Kirkpatrick et al., 'Overcoming catastrophic forgetting in neural networks', Proceedings of the National Academy of Sciences of the United States of America, vol. 114, no. 13, pp. 3521–3526, Mar. 2017, doi: 10.1073/PNAS.1611835114.

[24] R. Kemker, M. McClure, A. Abitino, T. L. Hayes, and C. Kanan, 'Measuring catastrophic forgetting in neural networks', in 32nd AAAI Conference on Artificial Intelligence, AAAI 2018, 2018. doi: 10.1609/aaai.v32i1.11651.

[25] T. L. Hayes, K. Kafle, R. Shrestha, M. Acharya, and C. Kanan, 'REMIND Your Neural Network to Prevent Catastrophic Forgetting', Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), vol. 12353 LNCS, pp. 466–483, 2020, doi: 10.1007/978-3-030-58598-3_28.

[26] C. Greco, B. Plank, R. Fernández, and R. Bernardi, 'Measuring Catastrophic Forgetting in Visual Question Answering', in Lecture Notes in Electrical Engineering, vol. 714, 2021. doi: 10.1007/978-981-15-9323-9_35.

[27] H. V. Hasselt, A. Guez, and D. Silver, 'Deep Reinforcement Learning with Double Q-Learning', in Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence (AAAI-16), 2016, pp. 2094–2100.

[28] Z. Wang, T. Schaul, M. Hessel, H. Van Hasselt, M. Lanctot, and N. De Freitas, 'Dueling Network Architectures for Deep Reinforcement Learning', 33rd International Conference on Machine Learning, ICML 2016, vol. 4, no. 9, pp. 2939–2947, 2016.

[29] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, 'Prioritized experience replay', 4th International Conference on Learning Representations, ICLR 2016 - Conference Track Proceedings, pp. 1–21, 2016.

[30] Josh Achiam, 'Deep Deterministic Policy Gradient', OpenAI Spinning Up. [Online]. Available: https://spinningup.openai.com/en/latest/algorithms/ddpg.t.

[31] P. Schulman, John and Levine, Sergey and Abbeel, Pieter and Jordan, Michael and Moritz, 'Trust region policy optimization', in International conference on machine, 2015, pp. 1889–1897. doi: 10.3917/rai.067.0031.

[32] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, 'Proximal Policy Optimization Algorithms', arXiv preprint arXiv:1707.06347, pp. 1–12, 2017.

[33] Z. Lin, T. Zhao, G. Yang, and L. Zhang, 'Episodic memory deep q-networks', in IJCAI International Joint Conference on Artificial Intelligence, 2018, pp. 2433–2439. doi: 10.24963/ijcai.2018/337.

[34] V. Mnih, M. Mirza, A. Graves, T. Harley, T. P. Lillicrap, and D. Silver, 'Asynchronous Methods for Deep Reinforcement Learning', in 33rd International Conference on Machine Learning, ICML 2016, 2016, pp. 1928--1937.

[35] W. Ko and D. E. Chang, 'A dual memory structure for efficient use of replay memory in deep reinforcement learning', in 2019 19th International Conference on Control, Automation and Systems (ICCAS), 2019, pp. 1483--1486.

[36] D. Zha, K. H. Lai, K. Zhou, and X. Hu, 'Experience replay optimization', IJCAI International Joint Conference on Artificial Intelligence, vol. 2019-Augus, pp. 4243–4249, 2019, doi: 10.24963/ijcai.2019/589.

[37] S. Zhang and R. S. Sutton, 'A deeper look at experience replay', arXiv preprint arXiv:1712.01275, 2017.

[38] P. Sun, W. Zhou, and H. Li, 'Attentive Experience Replay', in Proceedings of the AAAI Conference on Artificial Intelligence, 2020, pp. 5900–5907. doi: 10.1609/aaai.v34i04.6049.

[39] M. Brittain, J. Bertram, X. Yang, and P. Wei, 'Prioritized sequence experience replay', arXiv, 2019.

[40] T. de Bruin, J. Kober, K. Tuyls, and R. Babuska, 'The importance of experience replay database composition in deep reinforcement learning', Deep Reinforcement Learning Workshop, Advances in Neural Information Processing Systems (NIPS), pp. 1–9, 2015.

[41] B. Mavrin, H. Yao, and L. Kong, 'Deep Reinforcement Learning with Decorrelation', Mar. 2019, [Online]. Available: http://arxiv.org/abs/1903.07765.

[42] W. Olin-Ammentorp, Y. Sokolov, and M. Bazhenov, 'A Dual-Memory Architecture for Reinforcement Learning on Neuromorphic Platforms', Neuromorphic Computing and Engineering, vol. 1, p. 024003, 2021.

[43] N. Kamra, U. Gupta, and Y. Liu, 'Deep Generative Dual Memory Network for Continual Learning'. 2017. [Online]. Available: http://arxiv.org/abs/1710.10368.

[44] R. Liu and J. Zou, 'The Effects of Memory Replay in Reinforcement Learning', 2018 56th Annual Allerton Conference on Communication, Control, and Computing, Allerton 2018, pp. 478–485, 2019, doi: 10.1109/ALLERTON.2018.8636075.

[45] H. R. Maei, 'Convergent Actor-Critic Algorithms Under Off-Policy Training and Function Approximation'. arXiv, Feb. 21, 2018. Accessed: Nov. 24, 2023. [Online]. Available: http://arxiv.org/abs/1802.07842.

[46] V. R. Konda and J. N. Tsitsiklis, 'Actor-critic algorithms', Advances in Neural Information Processing Systems, pp. 1008–1014, 2000.

[47] J. Tarbouriech, M. Pirotta, M. Valko, and A. Lazaric, 'A Provably Efficient Sample Collection Strategy for Reinforcement Learning', pp. 1–33, 2020.

[48] G. Brockman et al., 'OpenAI Gym', arXiv preprint arXiv:1606.01540, pp. 1–4, 2016.

[49] P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger, 'Deep reinforcement learning that matters', in 32nd AAAI Conference on Artificial Intelligence, AAAI 2018, 2018.

[50] N. Raffin, Antonin and Hill, Ashley and Ernestus, Maximilian and Gleave, Adam and Kanervisto, Anssi and Dormann, 'Stable Baselines', GitHub repository, 2019, [Online]. Available: https://github.com/hill-a/stable-baselines.

[51] E. F. Morales and H. J. Escalante, 'A brief introduction to supervised, unsupervised, and reinforcement learning', in Biosignal processing and classification using computational learning and intelligence, Elsevier, 2022, pp. 111–129. Accessed: Nov. 25, 2023. [Online]. Available: https://www.sciencedirect.com/science/article/pii/B97801282012510001 78.

[52] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, and J. Pineau, An introduction to deep reinforcement learning, vol. 11. 2018. doi: 10.1561/2200000071.

# Real Time FPGA Implementation of a High Speed for Video Encryption and Decryption System with High Level Synthesis Tools

Ahmed Alhomoud

Department of Computer Sciences, Faculty of Computing and Information Technology,
Northern Border University, Rafha 91911, Saudi Arabia

*Abstract*—**The development of communication networks has made information security more important than ever for both transmission and storage. Since the majority of networks involve images, image security is becoming a difficult challenge. In order to provide real-time image encryption and decryption, this study suggests an FPGA implementation of a video cryptosystem that has been well-optimized based on high level synthesis. The MATLAB HDL coder and Vivado Tools from Xilinx are used in the design, implementation, and validation of the algorithm on the Xilinx Zynq FPGA platform. Low resource consumption and pipeline processing are well-suited to the hardware architecture. For real-time applications involving secret picture encryption and decryption, the suggested hardware approach is widely utilized. This study suggests an implementation of the encryption-decryption system that is both very efficient and area-optimized. A unique high-level synthesis (HLS) design technique based on application-specific bit widths for intermediate data nodes was used to realize the proposed implementation. For HLS, MATLAB HDL coder was used to generate register transfer level RTL design. Using Vivado software, the RTL design was implemented on the Xilinx ZedBoard, and its functioning was tested in real time using an input video stream. The results produced are faster and more area- efficient (target FPGA has fewer gates than before) than those of earlier solutions for the same target board.**

*Keywords*—*Security; encryption; decryption; AES; HDL coder; high levelsynthesis; FPGA; Zynq7000*

## I. INTRODUCTION

The development of strong, computationally light, and efficient encryption algorithms is therefore necessary to support the ongoing increase in data volume and throughput in Internet of Things applications, as well as video streaming, real-time video processing, mobile transmissions, and other related activities. This is because network security has become. A constant topic for business activities due to the advancements of information technology (IT) applications involving sensitive data [1,2]. For governments, banks, and high-security systems worldwide, the Advanced Encryption Standard (AES) continues to be the recommended encryption standard [3,4]. The last encryption technique is the most often used; it is utilized in 5G systems, WiMAX, Gigabit Ethernet, and Worldwide Interoperability for Microwave Access [5,6,7]. Furthermore, this algorithm can be effectively implemented on both software and hardware platforms; AES software versions give poorer physical security but require fewer resources [8, 9].

However, the increasing need for secure data transmissions at high speeds and volumes while maintaining physical security makes hardware implementation of the AES algorithm imperative [10,11]. The primary challenge with applications utilized in these domains is to ensure real-time system operation [2]. Implementing real-time functionality on a general-purpose computer is often unfeasible due to the inherent limitations of memory, CPU, and peripheral devices. Typically, numerous actions are executed on every pixel in the majority of image processing programs. The sequential execution of these operations by general-purpose processors has detrimental effects on both resource use and performance [2,3]. Nevertheless, FPGAs (Field Programmable Gate Arrays) possess the potential to function in a parallel manner in relation to hardware, setting them apart from conventional CPUs. Operations in FPGAs are partitioned into segments, allowing for concurrent execution of many operations [12,13]. Fig. 1 presents a complete system for real-time image and video processing using an embedded system.

The landscape of smart video encryption is evolving rapidly, demanding sophisticated analysis of live video streams to accurately identify objects, scenes, and critical events. This has led to the integration of advanced analysis mechanisms into the traditional image processing pipelines of modern video cameras [14]. However, the stringent constraints of real-time processing and low power consumption inherent in camera modules limit the complexity and number of operations that can be feasibly implemented [15]. To address this challenge, researchers have focused on a select few pre-processing tasks like motion estimation, image segmentation, and robust video encryption [15].

Recognizing the growing complexity of computer vision systems, designers are increasingly turning to higher-level programming and synthesis tools to expedite the development process and overcome hardware limitations. Two prominent tools in this arena are Simulink Hardware Description Language (HDL) Coder and Xilinx High-Level Synthesis (HLS) [16, 17]. Xilinx HLS stands out for its exceptional suitability for designing large-scale computer vision systems. It boasts the ability to seamlessly integrate pre-existing standard algorithms and offers comprehensive functional verification, ensuring accuracy and efficiency [17]. Additionally, HDL Coder empowers rapid synthesis and verification of diverse image processing methods, ranging from picture statistics gathering and custom filtering to color space conversion [16].

Fig. 1. Embedded real time image and video processing system.

Nevertheless, there is no specific support for picture segmentation tasks in the present toolbox version. To accomplish this, a Simulink model is created that increases this toolbox's capacity and supports this essential feature. While academics have recently proposed a number of advanced methods for image encryption, implement" encryption decryption algorithm" implemented [14, 15] into our proposed hardware in order to minimize the use of logic resources. Moreover, it has been shown that switching from moving averaging to weighted averaging reduces the amount of logic resources needed without sacrificing the accuracy of the results. As a result, the following can be used to summarize the contributions of the completed work: Creation of a synthesizable Simulink model for the AES-based cryptosystem, which isn't yet accessible as an intrinsic block in the Simulink HDL Coder/Vision HDL Coder toolbox (MATLAB R2018b). By substituting the weighted average for the moving average, which necessitates an expensive division operation, logic resource conservation is achieved. In comparison to earlier methods, this work presents a real-time implementation of video encryption and decryption on a Xilinx ZedBoard and shows that it is faster and uses less space on the FPGA. A unique high-level design technique was used to create the design, which synthesizes the design with intermediate signal widths limited by the application (input stimulus). The rest of this essay is structured as follows: An introduction to high- level synthesis is given in Section I. Related work is given in Section II. FPGA High level synthesis is given in Section III. The AES Encryption and Decryption architecture is explained in Section IV. Image and video acceleration is given in Section V. RTL design implementation mentioned in Section VI. Finally Section VII concludes the paper.

## II. RELATED WORK

Recently, numerous researchers have undertaken investigations on cryptographic algorithms inside the realm of Internet of Things (IoT). Reference [18] introduced a dynamic pass- word access approach for uniformly storing the key

matrix on all nodes. The sender did not need to transmit a basic key, but rather the storage coordinates of the key matrix. The receiver then extracted the key from the matrix based on these coordinates, therefore enhancing the security of key transmission. Reference [19] introduced a Very Large-Scale Integration (VLSI) design that incorporates a 64-bit data path for the lightweight cipher present. This architecture achieves excellent performance, low power consumption, and a compact footprint on FPGA, resulting in a high throughput rate. In order to fulfill the security demands of various application contexts, distinct techniques are necessary for the encryption and decryption system to handle data. Reference [20] developed a customizable encryption system that allowed users to choose their preferred encryption method from a range of options specified in the FPGA configuration file, hence enhancing the flexibility of the system. Reference [21] suggested a hybrid protocol architecture for Short Message Services (SMS) that incorporates AES and Rivest Cipher 4 (RC4) algorithms to enhance the security of smart houses. This solution offered secure communication in the IoT context, ensuring confidentiality and randomness. However, it incurred a certain level of cost in terms of both time and space. Hossain et al. In study [21] author developed a flexible encryption method on the FPGA. Users have the option to choose between the AES, Data Encryption Standard (DES), and 3DES algorithms for encryption, based on their specific needs. This design enhanced system adaptability, but incurred wastage of logic resources. The advent of dynamically reconfigurable technology offered a superior option to achieve a balance between system flexibility and hardware resource usage.

Many studies are devoted to creating specialized hardware accelerators that can be utilized to carry out specific tasks in applications related to image and video processing. To showcase the system's functionality, spatial filters were implemented on the embedded platform Zedboard [22]. The paper examines a recent study on a hardware accelerator for video processing, which was developed on an Altera Cyclone IV FPGA. The accelerator is engineered to possess reduced

processing and memory bandwidth demands. The research findings are detailed in citation [23, 24]. Platforms Based Design (PBD) analyzes the Virtex-5 FPGA's performance in real time while processing images and videos by looking for commonalities and differences among different design criteria. To implement an effective architectural solution, the Xilinx ML-507 platform runs a PBD. This system is capable of real-time 60 fps video capturing in VGA resolution [5]. On the Xilinx Zynq7000 System on a Chip (SoC), different designs have been constructed using the Histogram of Oriented Gradients, or HOG, method. These designs can process images with a resolution of 1920 x 1080 pixels and achieve a frame rate of 39.6 frames per second: [6]. The authors demonstrated how to use the OpenCV function on an ARM processor for hardware implementation. The several classifications of hardware accelerated systems that employ Field-Programmable Gate Arrays (FPGAs) for image analysis are concisely listed in reference [2]. FPGAs have the capability to perform concurrent execution of several threads, allowing them to effectively handle a diverse array of applications, including those commonly encountered in the automotive industry. When it comes to driver assistance systems (DA), most researchers have highlighted the improvements in image and video processing [14]. Claus et al. [15] proposed the utilization of dynamic partial reconfiguration (DPR) in driver assistance (DA) systems, employing Multiple Processor System-on-Chip (MPSoC) architecture, to enhance security in various driving situations. The "Autovision" architecture was designed to demonstrate the benefits of DPR by utilizing hardware accelerated engines. Several observations focus on the development of hardware-based real-time lane identification for advanced driving assistance systems. Using the Hough Transform, real-time lane detection is implemented at a clock speed of 100 MHz on the Xilinx Zynq-7000 platform. Thanks to its $480 \times 270$-pixels resolution, it can run at 130 fps per second. The implementation is performed via the Vivado HLS tool [16]. FPGAs are well-suited for complex video and image processing workloads, such as K-means clustering.

The process of dividing an image into segments and compressing it without any loss of data is referred to as picture segmentation and compressing without loss [17]. Edge detection is an advanced image processing technique mostly used in surveillance systems. The article authored by Kowalczyk et al. [8] examines the practical execution of 4K streaming of videos on Xilinx devices. A high-definition video streaming system utilizing quick prototyping has been developed, employing an FPGA-based edge detection design [9]. The research presents a comparative analysis and investigation into edge recognition filters implemented on Field-Programmable Gate Arrays (FPGAs) for real-time processing of video and images techniques [10]. Babu et al. [11] present a succinct analysis of the various classifications of FPGA architecture and their corresponding applications.

### III. FPGA High-Level Synthesis for Image Processing System with Matlab HDL Coder

After the text edit has been completed, the paper is ready for the template. Duplicate the template file by using the Save As command and use the naming convention prescribed by your conference for the name of your paper. In this newly created file, highlight all of the contents and import your prepared text file. You are now ready to style your paper; use the scroll down window on the left of the MS Word Formatting toolbar.

Fig. 2 displays the simplified block diagram of the proposed system architecture. The design is suitable for both image and video encryption and decryption, as it allows for the sequential streaming of information for each pixel. However, when it comes to video processing, the method may require a single pixel or several pixels.

The system comprises a range of video processing processes that can be controlled by specific algorithms, resulting in the faster execution of certain filters on the integrated blocks of the Xilinx SoC platform. The output is directly transferred to a video processing component and display component. In addition, the processed output is produced and displayed outside on an HDMI monitor. Xilinx implements various video processing accelerators in the programmable logic (PL) based on the design.

Fig. 2 shows the design schematic for the proposed system. This method is applicable to video and image processing due to the sequential processing of the input pixels. When encrypting and decrypting embedded videos, mega pixels are required. Built on the embedded Xilinx SoC platform blocks, it can run complex algorithms faster thanks to a video block that can be configured with unique algorithms. The video processing block and display component directly accept the output from the video and picture systems [16]. The Xilinx Zynq 7000 platform systems are utilized for the implementation of several FPGA video and image processing accelerators [5,14]. The Advanced Extensible Interface (AXI) in the architecture shown in Fig. 1 connects the processor responsible for transferring the incoming video to the video processing pipeline system. The ARM CPU and the USB camera communicate over interfaces and exchange data. The memory controller contains the frame data. The HDMI display exhibits the method when the video IP core has finished executing it on the frame. The pipeline is being maintained for future versions [10].

The necessity to create intricate DSP systems that call for specialized arithmetic units, such the addition-compare-select unit for the Viterbi decoder, gave rise to the MBD approach for FPGAs. A more refined degree of FPGA-based circuit optimization is needed for these computing units [17]. This degree of optimization is typically linked to conventional digital design systems for processing images and videos. In essence, model-based design describes the real-time interactions that a system will have with the analog environment. The bulk of these apps make use of the widely used Unified Modeling Language (UML). The methods by which these tools establish and describe a system differ. These tools may employ various implementation tactics, some of which may prove to be less effective than others. However, they assure rapid system prototyping, ensuring punctuality. Some factors that affect the choice of a tool are its level of flexibility, its availability as pre-built libraries and blocks, and its general understanding.

Fig. 2. Crypto-video system on FPGA based on Zynq7000 platform.

Matlab, Simulink Realtime Workshop, Arti-Real Time Studio, and Rhapsody from Ilogix are among the UML-based tools [15]. These instruments are used in the creation of embedded multiprocessor environments. There are two types of tools that facilitate the creation of HDL code for FPGA: block-based tools and C language-based tools that utilize blocks to produce HDL code from block diagrams. Following that, the HDL code is utilized by the hardware of the synthesis tool in order to put the system design into action on the FPGA computer. Synplify DSP, Xilinx System Generator, DSP Builder Altera, and Simulink HDL Coder are some of the tools that are predicated on the Simulink and MATLAB environments. The majority of these tools are based on these environments. These technologies ensure a sophisticated modeling environment for signal processing algorithms. Combining the IP cores from FPGA manufacturers with blocks from the Simulink library results in the creation of HDL code that is singular to the platform in question. The designer benefits from more freedom through the utilization of tools such as Simulink HDL Coder, which seamlessly incorporates MATLAB functions and m-block files. The designer generates a Simulink model and subsequently transforms it into the FPGA environment with these tools.

The second group of MBD tools uses C to construct a system design abstraction. Among these are the Handel-C from Celoxica and the Catapult C from Mentor Graphics. A popular tool for developing embedded systems on FPGA, the Simulink HDL coder is the main engine behind these products [9].

MathWorks introduced its hardware/software workflow for Zynq-7000, with a specific focus on Model-Based Design (MBD), in September 2013, as stated in references [10, 11]. According to the depicted process in Fig. 3, Simulink is utilized with HDL toolkit to create models that may effectively demonstrate a fully dynamic system. These encompass a Simulink model designed for algorithms specifically tailored for the Xilinx Zynq SoC platform. Additionally, it enables the rapid creation of software-hardware implementations for the Zynq platform directly from the algorithm and system architecture.



Fig. 3. Model based design prototyping with MATLAB / HDL coder.

The development of the suggested real-time video processing system is bifurcated into two components: 1) Designing the architecture of a video processing system, and 2) Designing algorithms for video processing. The initial section examines the primary elements that contribute to the video processing system on the Zynq platform, specifically focusing on the AXI4 Interfaces utilized for efficient data transmissions. The subsequent section delves into the strategies and enhancements implemented for constructing video processing algorithms using Vivado HLS [25]. The proposed approach is founded upon the following fundamental principle:

- The utilization of Simulink simulation by system designers and algorithm developers serves two purposes.

The designer's task is building models for an entire system, encompassing communications, image, and video processing components. The second purpose is to enable the division of the model into hardware and software components and achieve a favorable balance for high-level synthesis.

- The Xilinx Zynq 7000 platform can be equipped with high-speed I/O cores and IP cores through the utilization of HDL code generation from HDL coder TM.

- The Zynq Cortex-A9 cores can be programmed using the embedded coder from Simulink, which facilitates quick iteration of embedded software development.

- The ARM processor system and programmable logic with support for Xilinx Zynq 7000 may generate automatic AXI4 interface cores.

- Integration with subsequent processes, such as software compilation, generating the executable for the ARM, and creating bit streams using Xilinx implementation tools like Vivado, allows for a fast prototype process. These jobs can be directly downloaded to Zynq 7000 platform boards.

## IV. FUNDAMENTALS OF THE AES-128 ENCRYPTION / DECRYPTION ALGORITHM

Similar to the Data Encryption Standard (DES), the AES algorithm presented in Fig. 4 operates as a block cipher at the bit level. Each block length is set at 128 bits, whereas the key length can be any value between 128 and 256 bits [20]. Every 128-bit data block is divided into 16 bytes, which are then mapped onto a $4 \times 4$ array called state. Every byte in the state represents a 2 x 8 cardinality Galois Field (GF) element. The algorithm consists of n rounds, or iterations, depending on the length of the key. For a key length of 128 bits, 192 bits, or 256 bits, the number of rounds is 10, 12, or 14. With the exception of the final round, each encryption round consists of the following four operations:

- Substitute Bytes

- Shift Rows

- Mix Columns

- Add Round Key

Every operation is executed in turn throughout each round, with the exception of the first Add Round Key; the Mix Columns action is skipped in the last round (see Fig. 5).



Fig. 5. AES 128 encryption algorithm.

The Substitute Bytes step is a non-linear transformation in which the relationship between the key and the cipher-text is hidden by replacing each byte in the state array with the entry of a fixed 8-bit Substitution Box (Sbox), which is implemented as a lookup table with 2 8 words of 8 bits each. To prevent assaults based on basic algebraic features, the Sbox utilized is constructed from the multiplicative inverse over GF(28) and paired with an invertible geometric transformation, yielding a $16 \times 16$ bytes table (see Fig. 5). Based on the most and least significant nibbles of the 8-bit input data, the permutation is obtained T The bytes in each row are circularly shifted by a specified offset during the Shift Rows step's operation on the state array's rows. Every byte in the second row is moved one position to the left; similarly, the third and fourth rows are shifted by two and three bytes to the left, respectively. The first row remains unmodified. In the Mix Columns phase, each column of the state array is mixed linearly by treating it as a polynomial over GF(28). Each column is then multiplied, modulo z4 + 1, by a fixed polynomial (c(z) = 03z3 + 01z2 + 01z + 02). The relationship between the plain text and the ciphertext must be concealed using both the Mix Columns and Shift Rows methods.

## V. IMAGE AND VIDEO ACCELERATION WITH HIGH LEVEL SYNTHESIS

The functional diagram of the suggested system design is displayed in Fig. 6. Since the input pixels are processed sequentially, the approach can be used to both image and video processing. Mega pixels are needed in the embedded video processing process. It is constructed from a video block that can be programmed with customized algorithms, speeding up complicated and resource- and time-consuming algorithms on the embedded Xilinx SoC platform blocks. A video processing block and a display component receive the output directly from the video and image systems [16]. The Xilinx Zynq 7000 platform systems are used to implement the several FPGA video and image processing accelerators [5] One MBD tool that makes system modeling, analysis, and simulation possible



Fig. 4. AES encryption and decryption process.

is Simulink HDL coder. It provides the designer with an organized graphical environment that enables the creation of highly complicated system designs. Moreover, the user of this application can generate adaptable custom blocks using MATLAB functions. utilizing Simulink blocks, the designer can produce bit-precise and synthesizable HDL code from the model by utilizing the Simulink HDL coder. Altera Quartus II,

Synplify, and Xilinx Vivado are some of the tools that can be used to synthesis and map the obtained HDL code to the target FPGA board. There are numerous built libraries in the Simulink HDL coder [11]. Adders, multipliers, accumulators, integrators, multi-port switches, lookup tables, etc. are a few examples of these preset libraries.



Fig. 6. FPGA crypto system design.

A typical MBD design flow for implementing FPGA-based video and image processing system is shown in Fig. 6. this system is divided on two part,First part is to generate a complete video processing system based on Matlab/HDL Coder . This generated vivado project system is without encryption and decrytion IPs.The second part is designing with VHDL language the top level for encryption and decrytion block.After Encryption and decryption IPs simulation and verification. Theese IPs are added to the full complete vivado project for real time video process-ing.Finally the final combined system for video encryption and decryption is implemented and verified on zynq7000 paltform.

## VI. RTL DESIGN IMPLEMENTATION

### A. Simulation

To validate the accuracy of the encryptions and decryption implementation, a testbench has been created to compare two distinct ciphertexts generated by this implementation with the expected true ciphertexts provided in reference. The implementation successfully passes the verification process, and a snapshot of the waveform produced from the simulation using the vivado 2017.4 simulator is shown in Fig. 7.

### B. Synthesis and FPGA Implementation

In the third and final phase (see Fig. 8), the Simulink HDL Coder transforms the Software-Hardware model of the video processing system into an IP core that is compliant with the AXI4 streaming bus and is in the form of HDL (VHDL) source

code. Encryption and decryption IP designed with VHDL language are integrated using Vivado to the full System to generate the RTL video encryption and decryption design as presented in Fig. 8.

In order to confirm the functionality of this IP core in a real-time, practical setting, a Hard- ware–Software co-design (HW–SW) has been directly implemented on a 170MHz Xilinx Zynq-7000 AP SoC XC7Z020- CLG484 FPGA. A Full Vivado project is generated for the HW-SW by the Simulink HDL coder. The Xilinx Vivado tool (version 2017.4), along with all the hardware and software add-ons, can be used to implement this project. Additionally, a single bus connects the Sobel core and the Color transform IP core. The video processing system Vivado project now includes the AES encryption-decryption based on VHDL as an IP. The RTL design for our crypto-video system is shown in Fig. 8.



Fig. 7. RTL simulation for AES encryption.

Fig. 8. RTL video encryption and decryption design.

## C. Resource Utilization

The proposed system used the Zynq-7000 SOC to implement the cores of our video processing system. Table I lists the materials that were utilized. BRAM is used to store data values, firmware, and instruction memory.

The architecture of the processor, which includes control signals, internal registers, and mi- crocode, is defined by the consumed LUTs and DFFs. Our suggested system core operates best at about 296MHz with a throughput of 95FPS. The Zynq7000 SOC's resources are impacted by the intricacy of the applied design. The suggested architecture operates with a 1080-1920 frame input resolution. These findings support the Vivado-based reconfigurable SOC platform.Table I provides a breakdown of the resources used by the system, including Block RAM (BRAM) for data storage and Look-Up Tables (LUTs) and Flip-Flops (DFFs) for the processor architecture. This information is crucial for understanding the resource footprint of the design and its feasibility for different Zynq-7020 models.

## D. Encryption / Decryption Time Test

Encryption and dcreyption speed is crucial. The timer is implemented to precisely track encryption and decryption cycles, determining the exact number of operations needed for each process. This granular data allows us to optimize performance and ensure efficient data protection.

$$T_{Proc} = N_C \times T_{cycle} = \frac{N_c}{FPGA_{Freq}}$$

where, Nc, the total encryption cycles for the image, is calculated by multiplying the clock cycle duration (Tcycle = 1/FPGAFreq) by the total execution time.

For a 512x512 image, the proposed encryption decryption algorithm required approximately 53.96 million cycles (75ms) and decryption required 54.85 million cycles (80ms). This compares favorably to other solutions in the literature (see Table II).

TABLE I. VIDEO ENCRYTION AND DECRYPTION RESOURCES UTILISATION

| Xilinx Platform | Zynq 7000 XC7Z020-1CLG484C |
|---|---|
| Maximum Frequency | 296.789MHz |
| LUT-FF Pairs | 1104 |
| LUTs as Logic | 1104 |
| LUTs as Memory | 18 |
| Slice Registers | 264 |
| RAM 36/18 | 1 |
| Frame Rate | 95FPS |

TABLE II. PROPOSED SYSTEM TIME COMPARISON

| IP | Time (s) |
|---|---|
| Encryption Time | 0,075 |
| Decryption time | 0,08 |
| Image preprocessing time | 0.035 |
| **Total** | **0.19** |

## VII. CONCLUSION

Accurate implementation PPA details can be difficult, if not impossible, to obtain during the algorithm development phase in a normal design flow because doing so means putting a lot of work on the implementation team to complete experimental implementations. Algorithm developers employ imprecise estimations for PPA prediction because of this difficult task. The PPA objectives are frequently not met as a result of algorithm developers' inaccurate estimations. Algorithm developers an quickly obtain precise PPA information using the MATLAB connection with Stratus HLS, which enables measurement- driven algorithm improvement. The solution automates a large portion of the manual process from MATLAB to implementation details with minimal disruption to the current design flow. Additionally, this connection

automates the creation of more optimal RTL and micro-architectural exploration, resulting in shorter design deadlines and better PPA. According to preliminary findings, this process will essentially become the norm for creating and executing silicon-targeted algorithms. In this study, the application of typical AES encryption is investigated on a Xilinx ZedBoard equipped with a Zynq- 7000 SoPC. This effort concentrated on the encryption side of AES128; however, it would not be difficult to construct and test the decryption side as well. Xilinx Vivado High Level Synthesis was used to implement the AES after it was first coded in a high-level language. It will be easy to swiftly implement our design and perform changes that significantly raised the AES algorithm's throughput thanks to the Xilinx HLS tool. Additionally, HLS has the ability to enable thorough examination of a design's resource utilization in comparison to high-level code placement, as well as hardware testing during the early phases of design.future work will explore expanding the HLS-MATLAB integration to encompass the decryption side of AES as well as investigating its application to a wider range of algorithms across diverse domains. This continued exploration holds immense promise for revolutionizing the way the proposed system is developed and deployed efficiently, high-performance hardware solutions.

### REFERENCES

[1] Li, L.; Li, S. High throughput AES encryption/decryption with efficient reordering and merging techniques. In Proceedings of the 2017 27th International Conference on Field Programmable Logic and Applications (FPL), Gent, Belgium, 4–6 September 2017; IEEE: Piscataway, NJ, USA, 2017; pp. 1–4.

[2] Babu, P., and Parthasarathy, E. (2021). Hardware acceleration of image and video processing on Xilinx zynq platform. Intell. autom. soft comput, 30(3).

[3] Elsayed, G., Soleit, E., and Kayed, S. (2023). FPGA design and implementation for adaptive digital chaotic key generator. Bulletin of the National Research Centre, 47(1), 122.

[4] Li, K., Li, H., and Mund, G. (2023). A reconfigurable and compact subpipelined architecture for AES encryption and decryption. EURASIP Journal on Advances in Signal Processing, 2023(1), 1-21.

[5] Visconti, P.; Capoccia, S.; Venere, E.; Vela´zquez, R.; Fazio, R.d.10 Clock-Periods Pipelined Implementation of AES-128 Encryption- Decryption Algorithm up to 28 Gbit/s Real Throughput by Xilinx Zynq UltraScale+ MPSoC ZCU102 Platform. Electronics 2020, 9, 1665. https://doi.org/10.3390/electronics9101665.

[6] Abd El-Maksoud, A. J., Abd El-Kader, A. A., Hassan, B. G., Rihan.

[7] N. G., Tolba, M. F., Said, L. A., ... and Abu-Elyazeed, M. F. (2020). FPGA implementation of integer/fractional chaotic systems. Multimedia Security Using Chaotic Maps: Principles and Methodologies, 199-229.

[8] Wang, D., Lin, Y., Hu, J., Zhang, C., and Zhong, Q. (2023). FPGA Implementation for Elliptic Curve Cryptography Algorithm and Circuit with High Efficiency and Low Delay for IoT Applications. Micromachines, 14(5), 1037.

[9] Maazouz, M., Toubal, A., Bengherbia, B., Houhou, O., and Batel.

[10] N. (2022). FPGA implementation of a chaos-based image encryption algorithm. Journal of King Saud University-Computer and Information Sciences, 34(10), 9926-9941.

[11] Rajasekar, P.; Haridas, M. Efficient FPGA implementation of AES 128 bit for IEEE 802.16e mobile WiMax standards. Circuits Syst. 2016, 7, 371–380.

[12] Elsayed G, Kayed SI (2022) A comparative study between MATLAB HDL Coder and VHDL for FPGAs design and implementation. J Int Soc Sci Eng 4:92–98.

[13] Al-Musawi, W. A., Wali, W. A., and Al-Ibadi, M. A. (2021, July). Implementation of Chaotic System using FPGA. In 2021 6th Asia-Pacific Conference on Intelligent Robot Systems (ACIRS) (pp. 1-6). IEEE.

[14] Del-Valle-Soto, C.; Vela´zquez, R.; Valdivia, L.J.; Giannoccaro, N.I.; Visconti, P. An Energy Model Using Sleeping Algorithms for Wireless Sensor Networks under Proactive and Reactive Protocols: A Performance Evaluation. Energies 2020, 13, 3024.

[15] Noorbasha, F.; Divya, Y.; Poojitha, M.; Navya, K.; Bhavishya, A.; Rao, K.; Kishore, K. FPGA design and implementation of modified AES based encryption and decryption algorithm. Int. J. Innov. Technol. Explor. Eng. 2019, 8, 132–136.

[16] Ghodhbani, R., Saidani, T., Alhomoud, A., Alshammari, A., and Ahmed, R. (2023). Real Time FPGA Implementation of an Efficient High Speed Harris Corner Detection Algorithm Based on High-Level Synthesis. Engineering, Technology and Applied Science Research, 13(6), 12169-12174.

[17] Sikka, P., Asati, A. R., and Shekhar, C. (2021). Real time FPGA implementation of a high speed and area optimized Harris corner detection algorithm. Microprocessors and Microsystems, 80, 103514.

[18] Tsai, Y. H., Yan, Y. J., Hsiao, M. H., Yu, T. Y., and Ou-Yang, M. (2023). Real-Time Information Fusion System Implementation Based on ARM-Based FPGA. Applied Sciences, 13(14), 8497.

[19] Park, J.; Park, Y. Symmetric-Key Cryptographic Routine Detection in Anti-Reverse Engineered Binaries Using Hardware Tracing. Electronics 2020, 9, 957.

[20] Ghodhbani, R., Horrigue, L., Saidani, T., and Atri, M. (2020). Fast FPGA prototyping based real-time image and video processing with high- level synthesis. International Journal of Advanced Computer Science and Applications, 11(2).

[21] Bellemou, A.M.; Garc´ıa, A.; Castillo, E.; Benblidia, N.; Anane, M.; A´lvarez-Bermejo, J.A.; Parrilla, L. Efficient Implementation on Low Cost SoC-FPGAs of TLSv1.2 Protocol with ECCAES Support for Secure IoT Coordinators. Electronics 2019, 8, 1238.

[22] MathWorks, inc: HDL CoderTM User's GuideCOPYRIGHT 2012-2015 (2012) https://www.mathworks.com/help/hdlcoder/ Accessed 20 Feb 2023.

[23] Guerrieri, A., Upegui, A., and Gantel, L. (2023). Applications Enabled by FPGA-Based Technology. Electronics, 12(15), 3302.

[24] Sankar D, Lakshmi S, Babu C, Mathew K (2023) Rapid prototyping of predictive direct current control in a low-cost fpga using hdl coder. Int J Power Energy Syst 43(10):1–9. https://doi.org/10.2316/J.2023.203-0437.

[25] Yoon, I., Joung, H., and Lee, J. (2016). Zynq-based reconfigurable system for real-time edge detection of noisy video sequences. Journal of Sensors, 2016.

# A Method to Increase the Analysis Accuracy of Stock Market Valuation: A Case Study of the Nasdaq Index

Haixia Niu*

School of Finance, Henan Finance University; Zhengzhou Henan, 451464, China

*Abstract*—For a significant period, conventional methodologies have been employed to assess fundamental and technical aspects in forecasting and analyzing stock market performance. The precision and availability of stock market predictions have been enhanced by machine learning. Various machine learning methods have been utilized for stock market predictions. A novel, optimized machine-learning approach for financial market analysis is aimed to be introduced by this study. A unique method for improving the accuracy of stock price forecasting by incorporating support vector regression with the slime mould algorithm is presented in the present work. Other optimization algorithms were employed to enhance the prediction accuracy and the convergence speed of the network, which were Biogeography-based optimization and Gray Wolf Optimizer. An assessment of the proposed model's effectiveness in predicting stock prices was conducted through research employing Nasdaq index data extending from January 1, 2015, to June 29, 2023. Substantial improvements in accuracy for the proposed model were indicated by the results compared to other models, with an R-squared value of 0.991, a root mean absolute error of 149.248, a mean absolute percentage error of 0.930, and a mean absolute error of 116.260. Furthermore, not only is the prediction accuracy enhanced by the integration of the proposed model, but the model's adaptability to dynamic market conditions is also increased.

*Keywords—Machine learning; Nasdaq index; support vector regression; gray wolf optimizer; slime mould algorithm*

## I. INTRODUCTION

The stability and safety of the stock market are considered of utmost importance, as it is regarded as a crucial element within national economies [1] [2]. The behaviors and consequences of the stock markets have become a vital area of scholarly investigation owing to the potential risks associated with them [3]. The forecasting of stock price trajectories is deemed a crucial responsibility as it serves the dual purpose of maintaining stability within financial markets by regulators and enabling investors to make informed choices while mitigating potential risks. The utilization of uncertain prediction processes and the subsequent generation of erroneous predictions can potentially lead to significant hazards [4]. Hence, the development of a robust and compelling predictive model is deemed crucial for mitigating potential hazards. The issue of stock market uncertainty is addressed by other theories, while conventional forecasting techniques rely on patterns that exhibit consistent behaviour across time. The aforementioned approach fails to account for the inherent volatility of the stock market, and when combined with the multitude of variables involved, the task of predicting stock values becomes a complex endeavor. Nevertheless, the emergence of machine learning (ML) [5] [6] is being considered. A comprehensive approach that utilizes a variety of algorithms to optimize performance in various situations is demonstrated by the solution offered. This emerging advancement exhibits considerable promise in its capacity to fundamentally transform our methodologies for forecasting stock market trends. The notion that trustworthy information can be discerned and patterns within a given dataset can be identified through machine learning is commonly accepted [7].

The anticipation of stock market trends has been an enduring area of interest for scholars, and machine learning methodologies are progressively assuming a more prominent role in this regard. By conducting a comparative analysis of state-of-the-art machine learning methods using a decade of daily historical data from the top 10 equities on the Casablanca Stock Exchange, our study contributes to this ongoing dialogue. It is worth mentioning that ensemble learning has been utilized in the past for this purpose, as demonstrated by Bilal et al. [8]. These efforts employed various classifiers, including ridge regression, LASSO regression, support-vector machine (SVM), k-nearest neighbors, random forest, and adaptive boosting. SVM, adaptive boosting, random forests, and SVM were discovered to perform exceptionally well in short-term forecasting, demonstrating the effectiveness of ensemble learning across a range of prediction horizons. Expanding upon the aforementioned groundwork, Sonkavde et al. [9] investigated a variety of methodologies including time series analysis, ensemble algorithms, deep learning, and supervised and unsupervised machine learning, in order to address the complexities associated with stock price classification and prediction. Furthermore, a comprehensive analysis of the Nasdaq stock market was presented by Ashfaq et al. [10], who employed a variety of machine learning regressors to forecast the opening prices of specific companies the following day. For the purpose of evaluation, they utilized metrics such as the mean square error and coefficient of determination, thereby enhancing our comprehension of predictive modelling within this particular domain. Agrawal et al. [11] presented a seminal study in which they described an algorithmic approach for predicting the stock market that utilized deep learning and non-linear regression. The researchers' investigation, which utilized a decade's worth of data from the New York Stock Exchange and Tesla Stock Price, demonstrated that their proposed solution outperformed currently available machine learning algorithms. However, in their inability to adapt to the volatile and unpredictable nature of financial markets, these methodologies frequently encountered obstacles. The trade-offs that existed between accuracy and computational efficiency exposed this dilemma.

The advantages and disadvantages of these literatures are presented in the Table I. In order to fill this void, our research presents an enhanced machine-learning methodology that integrates the slime mould algorithm and support vector regression in a synergistic fashion. The objective of this innovative approach is to surmount the limitations identified in prior investigations through the improvement of forecast accuracy and flexibility in the face of market volatility. Our methodology signifies a substantial deviation from traditional approaches, enabling a more intricate examination of intricate market data and surmounting the constraints intrinsic in conventional models. Our objective is to provide a comprehensive resolution to the persistent difficulty of precisely forecasting stock market fluctuations within a dynamic financial environment.

TABLE I. THE ADVANTAGES AND DISADVANTAGES OF THE RELEVANT LITERATURE

| Authors | Advantages | Disadvantages |
|---|---|---|
| Bilal et al. [8] | Effective for short-term forecasting across a range of prediction horizons. | Lack of adaptability to unpredictable and volatile stock market. |
| Sonkavde et al. [9] | Addresses complexities of stock price classification and prediction. | Unable to adapt to stock markets that are volatile and unpredictable. |
| Ashfaq et al. [10] | Enhanced understanding of predictive modeling. | Unableness to adjust to unexpected and chaotic stock markets. |
| Agrawal et al. [11] | Outperformed existing machine learning algorithms. | Inability to adapt to volatile and unpredictable financial markets. |

Linear regression is one of the machine learning models most frequently utilized for forecasting. A valuable statistical technique that can be employed to make predictions about a numerical outcome is provided by linear regression [12]. Artificial neural networks (ANN) [13], Support vector machines (SVM) [13], Decision trees [14], Random Forest [15], Logistic regression [16], Gradient boosting [17], time series forecasting [18], and more. All of these models possess the limitations and advantages that exert a substantial influence on the accuracy of predictions.

The method of research employed in this study is support vector regression (SVR), which is a resilient form of supervised learning. In SVR, an effort is made to minimize both structural risk and empirical risk while also reducing the range of trust in training examples. A significant level of efficacy is demonstrated by the previously mentioned methodology in addressing intricate nonlinear issues, particularly those characterized by a limited sample size. Notable efficacy in addressing intricate nonlinear issues, particularly those characterized by a limited sample size, is demonstrated by the mentioned methodology. A crucial role is played by SVR in mitigating risk and enhancing the overall predictive accuracy of future samples. This, in turn, facilitates the extraction of useful insights and enhances the decision-making process [19]. An in-depth understanding of the principles and advantages of SVR is imperative for achieving ideal outcomes and accomplishments in several fields, such as machine learning, data science, and related domains.

The behavior and performance of a model during training are significantly influenced by a range of hyperparameters. Several factors, such as model complexity, regularization power, and learning rate, among other considerations, are encompassed by the hyperparameters. To optimize the model's efficacy, the careful selection and meticulous refinement of appropriate hyperparameters are imperative. Engaging in this practice has the potential to greatly affect the precision, resilience, and the applicability of the model, enabling it to more effectively align with the dataset on which it is being trained. Hence, the meticulous fine-tuning of the hyperparameters is an essential stage in achieving the utmost performance of the model. Various techniques and algorithms are employed for the optimization of hyperparameters for models, such as grey wolf optimization (GWO) [20] [21], Biogeography-based optimization (BBO) [22], Grasshopper optimization algorithm (GOA) [23], Slime mould algorithm (SMA) [24], Moth flame optimization (MFO) [25], and more. Some of the above optimization techniques are nature-inspired. Three strategies were applied to improve the proposed model's hyperparameters: Biogeography-based optimization, grey wolf optimization, and the Slime mould algorithm.

The GWO algorithm, as proposed by Mirjalili et al. [21], is inspired by the hierarchical structure of leadership and hunting patterns observed in grey wolves. Fundamentally, the grey wolf species are categorized into four distinct groupings, namely alpha, beta, delta, and omega. A methodology known as BBO, proposed by Nazari et al. [26], is utilized. A pioneering methodology called SMA was presented in a study conducted by Chen et al. This methodology is inspired by the behavioural patterns exhibited by slime mould organisms in natural environments. One intriguing attribute of slime mould is its capacity to perceive the olfactory cues associated with the presence of food particles in the surrounding environment. The olfactory perception of food odors in the surrounding environment is a notable characteristic that contributes to the intriguing nature of slime mold. The SMA approach aims to replicate this inherent mechanism to optimize its efficacy in attaining its goals.

In the present study, a comprehensive dataset covering the period from January 2015 to June 2023 was analyzed using many models. To ensure the accuracy and reliability of the produced outputs, thorough training was conducted for the SVR method, taking into account a wide array of input factors. The criteria included several factors such as daily transaction volume, high and low prices, as well as the opening and closing prices. To assess the accuracy of the model's results, a comprehensive testing procedure was conducted, employing the same parameters as those used during the initial phase. The outcome of this intensive training and evaluation procedure yields a model capable of furnishing traders and investors with invaluable market insights, facilitating well-informed decision-making, and ultimately fostering profitable investments. The variable data was acquired from the stock market of the Nasdaq. In order to achieve its objectives, the research employed a methodology comprising multiple analytical stages. The research paper presents a comprehensive examination of the data source and its relevant elements in the subsequent section. The data are analyzed utilizing a variety of

methods, including the SMA optimizer, evaluation metrics, and the SVR model. Following the presentation of the analyses' results in the third section, those obtained from alternative methodologies are contrasted. The study concludes with a concise presentation of the results in the concluding section.

## II. METHODOLOGY

### A. Support Vector Regression

A very effective technique utilized in the domain of machine learning, particularly for nonlinear classification. The SVR has become a widely adopted technique among data scientists due to its ability to handle inputs with high dimensions effectively [27] [28]. The SVR, as seen in Fig. 1, has demonstrated remarkable efficacy in resolving classification problems, prompting its extension to tackle challenges associated with regression. Similar to its predecessor, the present improvement, referred to as SVR, demonstrates exceptional proficiency in managing intricate data sets. SVR and its many extensions have emerged as vital tools in the field of machine learning due to their ability to provide a flexible and accurate methodology for addressing regression and classification problems [29].



Fig. 1. The diagram of the SVR.

The initial description pertains to a linear function, denoted as f(x), which is characterized by the following mathematical form:

$$f(x) = w \times x + b \qquad (1)$$

where, $x$ is the input vector, $b$ is a constant that has to be calculated, and $w$ is the vector holding the parameters. In instances involving nonlinear problems, the data is transformed into a higher-dimensional space by the utilization of a nonlinear kernel.

$$f(x) = w\phi \times (x) + b \qquad (2)$$

where, φ(x) is the kernel function.

The utilization of a feature space with larger dimensions can be employed to facilitate the mapping of data, hence enabling the implementation of a linear regression approach.

The coefficients $w$ and $b$ can be obtained by minimizing a given function as follows:

$$min\frac{1}{2} \parallel w \parallel^2 + c \sum_{i=1}^{N} (\xi_i + \xi_i^*) \qquad (3)$$

The exposure to:

$$\begin{cases} y_i - \langle w, x_i \rangle - b \leq \varepsilon + \xi_i \\ \langle w, x_i \rangle + b - y_i \leq \varepsilon + \xi_i^* \\ \text{with } \xi_i, \xi_i^* \geq 0, i = 1, \dots . N \end{cases} \qquad (4)$$

where, the positive and negative errors are denoted by $\xi_i$ and $\xi_i^*$, respectively. The constant C>0 is a hyperparameter that enables the adjustment of the trade-off between the permissible error and the flatness of the function $f(x)$.

### B. Biogeography-based Optimization

The BBO algorithm operates by simulating the movement of various species, which is determined by the suitability of their respective habitats. This process allows for a comprehensive analysis of the complex relationships between different species and their environments, ultimately leading to more refined and accurate predictions about ecological patterns. In the context of an optimization issue, it might be argued that a solution bears resemblance to a habitat. An optimal approach for addressing population concerns involves the establishment of densely populated habitats that offer superior living circumstances for various species compared to alternative habitats. The environment in which living animals encounter significant challenges is where the least optimal solution within the population is found. Therefore, the superior solutions can attract the inferior solutions due to their shared characteristics. The method of sharing features is accomplished by the utilization of the outlined operators.

Operation of Migration: Migration is a procedure through which a poorer habitat is swapped with a better one based on emigration rates. The quantification of the influx of species in a given area is referred to as the immigration rate. The migration incidence is anticipated to be greater under a more favourable solution compared to a less favourable option.

On the contrary, the quantitative assessment of the number of individuals within a species that depart from their environment is known as the rate of migration. Consequently, the emigration rate is expected to be greater under a suboptimal solution compared to an optimal option. The basic form of BBO has utilized straight paths as shown in Eq. (4).

$$\mu_k = \frac{E \times k}{n} \lambda_k = I \left(1 - \frac{k}{n}\right) \qquad (5)$$

$\boldsymbol{\mu_k}$: Migration amount of $\mathbf{k}^{th}$ habitat.

$\boldsymbol{\lambda_k}$: Migration amount of $\mathbf{k}^{th}$ habitat.

I: Maximum immigration amount.

E: Maximum emigration rate.

$\boldsymbol{n =}$: Maximum number of species that a habitat can support.

K: Number of species count.

Mutation: Mutation in BBO may be likened to an abrupt alteration in the environmental circumstances experienced by living organisms, such as those caused by natural disasters like earthquakes, volcanic eruptions, or tornadoes. Similarly, the random modifications in the genetic makeup of a species prompt its migration to a new habitat, as the previous habitat becomes unsuitable for its survival.

### C. Gray Wolf Optimizer

The GWO is a unique optimization approach that has been developed through the utilization of a meta-heuristic method. The strategy, initially introduced by Mirjalili et al. [18], emulates the social hierarchy and hunting strategies employed by grey wolves. Alpha is often regarded as the most optimum choice, whereas Omega is positioned as the final challenger inside the hierarchical framework of leadership.

Three principal hunting techniques are utilized by the method in order to emulate the behaviours of wolves: The action of following, confining, and assaulting prey. To simulate the locomotion patterns of grey wolves during hunting activities in their natural habitat, the following relate was employed:

$$\vec{D} = |\vec{C} \cdot \vec{X}_p(t) - \vec{X}(t)|$$

$$\vec{X}|(t+1) = \vec{X}_p(t) - \vec{A} \cdot \vec{D} \tag{6}$$

In which, $t$ is the current iteration, $\vec{D}$ denotes movement, $\vec{X}_p$ denotes prey location, $\vec{A}$ and $\vec{C}$ denotes coefficient vectors, and $\vec{X}$ denotes the grey wolf's position. The coefficient vectors ($\vec{A}$ and $\vec{C}$) are constructed using the relationships shown below.

$$\vec{A} = 2\vec{a} \cdot \vec{r}_1 - \vec{a}$$

$$\vec{C} = 2 \cdot \vec{r}_2 \tag{7}$$

The spatial allocation of novel search representatives pertaining to omegas is modified by using data derived from alpha, beta, and delta in the following manner:

$$\vec{D}_\alpha = |\vec{C}_1 \cdot \vec{X}_\alpha - \vec{X}|, \vec{D}_\beta = |\vec{C}_2 \cdot \vec{X}_\beta - \vec{X}|, \vec{D}_\delta = |\vec{C}_3 \cdot \vec{X}_s - \vec{X}| \tag{8}$$

$$\vec{X}_1 = \vec{X}_\alpha - \vec{A}_1 \cdot \vec{D}_u, \vec{X}_2 = \vec{X}_\beta - \vec{A}_2 \cdot \vec{D}_\beta, \vec{X}_3 = \vec{X}_b - \vec{A}_3 \cdot \vec{D}_\delta \tag{9}$$

$$\vec{X}(t+1) = \frac{\vec{X}_1 + \vec{X}_2 + \vec{X}_3}{3} \tag{10}$$

The wolves, denoted by the subscripts $\alpha, \beta,$ and $\delta$, converge to initiate a conclusive assault in order to accomplish the objective successfully. The variable $\vec{a}$ is used in order to replicate the previous assault by altering a value from 2 to 0. On the other hand, the variable a represents a random variable that falls between the range of $-2\vec{a}$ and $2\vec{a}$. Consequently, if $\vec{a}$ is lowered, it will also lead to a decrease in the value of $\vec{A}$. The wolves were compelled to grasp their prey due to the factor tightly $|\vec{A}| < 1$. Grey wolves engage in cooperative hunting strategies by forming packs and exhibiting a hierarchical social structure. These packs are led by an alpha wolf, who guides the group's activities. The pack members disperse to forage for food individually, and then afterwards regroup to launch coordinated attacks. Wolves may separate in search of prey when $|\vec{A}|$ has a random value greater than unity. The wolf count and generation number are considered to be the two most critical configuration parameters for the GWO algorithm. This means that the total number of objective function evaluations will be equal to the wolf population times the size of the generation or,

$$OFEs = N_W \times N_G \tag{11}$$

### D. Slime Mould Algorithm

Li et al. introduced a unique approach called SMA, which draws inspiration from the behaviour of slime mould seen in natural environments [24]. The slime mould, which is represented in Fig. 3, uses its olfactory perception and detects the volatilized scent of nourishment in the air in order to approach its prey. Fig. 2 provides a comprehensive illustration of the general characteristics of SMA.

The behaviour of the slime mould may be mathematically described by the equation that goes as follows:

$$\overrightarrow{X(t+1)} = \begin{cases} \overrightarrow{X_b(t)} + \vec{v_b} \cdot \left( \vec{W} \cdot \overrightarrow{X_A(t)} - \overrightarrow{X_B(t)} \right) & r < p \\ \vec{v_c} \cdot \overrightarrow{X(t)} & r \geq p \end{cases} \tag{12}$$

In which, $X_b(t)$ reflects the specific area of the slime mould that now exhibits the greatest level of odour, the variables $X(t)$ and $X(t+1)$ represent the locations of the slime mould in iteration $t$ and $t+1$, respectively. $X_A(t)$ and $X_B$ represent two randomly selected sites of slime mould. The variable $v_b$ varies over time within the range $[-a, a]$, where r is a random number between 0 and 1. The parameter $p$ is specified as $(a = \operatorname{arctanh}(-(\frac{t}{\max\_t}) + 1))$, and $v_c$ is a linearly decreasing parameter ranging from 0 to 1.

$$p = tanh |S(i) - DF| \quad i = 1, 2, \dots, n \tag{13}$$

where, $DF$ denotes the fittest iteration overall, and $S(i)$ denotes the fitness of $\vec{X}$. Following is a definition of the weight $W$ equation:

$$\overrightarrow{W(smell\ index(l))} =$$

$$\begin{cases} 1 + r \cdot log \left( \frac{bF - S(i)}{bF - wF} + 1 \right), condition \\ 1 - r \cdot log \left( \frac{bF - S(i)}{bF - wF} + 1 \right), others \end{cases} \tag{14}$$

$$smell\ index = sort(S) \tag{15}$$

The variable $S(i)$ represents the initial half of the population in the given equation. $bF$ for the best fitness, $wF$ for the poorest fitness, and the scent index for the sorted fitness values. By utilizing the formula provided, the spatial coordinates of the slime mould are revised:

$$\overrightarrow{X^*} = \begin{cases} rand(UB - LB) + LB & rand < z \\ \overrightarrow{X_b(t)} + \vec{v_b} \cdot \left( \vec{W} \cdot \overrightarrow{X_A(t)} - \overrightarrow{X_B(t)} \right) & r < p \\ \vec{v_c} \cdot \overrightarrow{X(t)} & r \geq p \end{cases} \tag{16}$$

Fig. 2.   The diagram of the Slime mould algorithm.

In this context, the parameter denoted as $z$ is constrained to a range between 0 and 0.1. The terms *LB* and *UB* refer to the lower and upper borders of the search interval correspondingly.



Fig. 3.   The illustration of the SMA algorithm.

*E.  Dataset Collection*

In order to conduct a thorough analysis, it is imperative to include trade volume and Open, High, Low, and Close (OHLC) prices during a designated period. Data for this investigation was collected via Yahoo Finance from January 2, 2015, to June 29, 2023. The Nasdaq, which was founded in 1971, is notable for being the inaugural electronic stock exchange in the globe. It distinguishes itself from conventional exchanges by functioning exclusively electronically and utilizing a digital infrastructure that optimizes the trading process. This allows for efficient and rapid transaction processing. In addition to technology firms, the Nasdaq Composite Index comprises companies from consumer services, healthcare, and a variety of other sectors. Due, in part, to its emphasis on emerging sectors and technology, the Nasdaq is a centre for innovative and expanding businesses, given its status as a major participant in international financial markets. Investors routinely track the Nasdaq Composite Index in order to assess the financial well-being of technology and growth firms, as well as to deduce wider economic patterns and investor sentiment.

## F. Dataset Description

The process of ensuring the high quality of the raw data is a crucial and fundamental stage in the pursuit of obtaining meaningful insights. The process of data preparation is of utmost importance in attaining this objective. The process encompasses a variety of tasks, such as the elimination of undesirable data, the establishment of standardized formats for widespread applicability, and the arrangement of information in a manner that promotes the retrieval of valuable insights. This has special significance in initiatives that entail substantial quantities of data because data quality takes precedence over mere numerical values. Data preparation activities encompass a range of tasks, such as encoding categorical data, cleaning and organizing data, scaling, standardization, and normalization in accordance with established industry standards. By engaging in these activities, the precision and dependability of the insights derived from the data can be enhanced. To achieve data scaling and normalization, Min-Max scalers were utilized in the project's data pre-processing step. This approach facilitated the removal of inconsistencies, as well as the measurement of null, missing, and unknown values. The methodology was illustrated using the data from the Nasdaq index. The dataset encompasses the timeframe spanning from January 2015 to June 2023 and has undergone several preparatory procedures, such as standardization.

The process of transforming numerical attributes inside a dataset to a specific range, commonly ranging from zero to one, is referred to as feature scaling. This technique is alternatively recognized as Min-Max normalization or data preparation. The objective is to maintain the relative relationships between the values while ensuring that all features are brought to a comparable scale. The consideration of input feature quantity holds particular significance in machine learning algorithms that exhibit sensitivity towards this aspect. The formula for the data normalization procedure is as follows:

$$\text{XScaled} = \frac{(X - Xmin)}{(Xmax - Xmin)} \qquad (17)$$

## G. Evaluation Metrics

Evaluation metrics are commonly used in the fields of statistics and machine learning to assess the effectiveness of forecasting models quantitatively. They aid in evaluating the predictive accuracy of a model on data that has not yet been observed. The selection of the optimal evaluation metric is contingent upon the particular analytical objectives and the nature of the predictive task at hand. Performance indicators such as R-squared ($R^2$), Root Mean Square Error (RMSE), Mean Absolute Percentage Error (MAPE), Mean Squared Error (MSE), Relative Squared Error (RSE), and Mean Absolute Error (MAE) were utilized in this study to evaluate the prediction accuracy of the developed forecasting models. Below, a compilation of mathematical formulas pertaining to these metrics is provided:

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \qquad (18)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}} \qquad (19)$$

$$MAPE = \left(\frac{1}{n}\sum_{i=1}^{n}\left|\frac{y_i - \hat{y}_i}{y_i}\right|\right) \times 100 \qquad (20)$$

$$MAE = \frac{\sum_{i=1}^{n}|y_i - \hat{y}_i|}{n} \qquad (21)$$

$$RSE = \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)}{\sum_{i=1}^{n}(\bar{y} - \hat{y}_i)} \qquad (22)$$

$$MSE = \frac{1}{N}\sum_{k=0}^{n}\binom{n}{k}(Fi - Yi)b^2 \qquad (23)$$

## III. RESULT AND DISCUSSION

### A. Statistical Results

Table II displays the statistical characteristics, such as mean, std., min, 25%, 50%, 75%, max, and variance. The central tendency of a data set is quantified by the mean, a statistical term. The calculation involves the sum of all values inside the dataset, followed by the division of the resulting sum by the total count of values. The dispersion or spread of data points relative to the mean is quantified by the standard deviation, a statistical metric. The minimum, minimum value, or minimal observation is the term used for the smallest data point or value in a dataset. It assists in finding the lowest number within a batch of data as well as understanding the range and distribution of the data. The degree to which individual data points in a dataset differ from the average (mean) of the dataset is described by the statistical concept of variance. A low variance denotes that data points are close to the mean, whereas a high variance suggests that data points are dispersed from the mean. Quantiles, specifically the 25th, 50th, and 75th percentiles, play a pivotal role in comprehending the data distribution. As indicated by the 25th percentile (also referred to as the first quartile), 25% of the observed values are situated below this threshold. As an illustration, a 25% percentile of 5776.33 for the opening price indicates that 25% of the dataset's initial prices fall below this threshold. As the 50th percentile, which is also equal to the median, divides the dataset in half, it is a significant indicator. With a median value of 7833.27, the proportion of closing prices that occur either above or below this value is half. The third quartile, or 75th percentile, indicates that 75% of the data are located below this value. A 75% percentile of the close price (1,590.78) indicates that 75% of the closing prices fall below this threshold. These quantiles offer valuable insights regarding the market's overall trend and stability. A significant disparity between the 25th and 75th percentiles, for instance, may suggest that stock prices are more volatile. A comprehension of the distribution of the majority of data points, namely stock prices, can assist analysts and investors in recognising customary price ranges, detecting anomalies, and identifying substantial changes in market trends.

The closing price data is shown in Fig. 4, whereby it has been partitioned into two distinct zones for the purposes of training and testing. This strategy ensures the precision of the data while aiding consumers in acquiring reliable insights.

TABLE II.        RESULTS OF THE STATISTICS FOR THE OHCLV MODELS THAT WERE PRESENTED

|          | Open      | High      | Low       | Volume   | Close     |
|----------|-----------|-----------|-----------|----------|-----------|
| mean     | 8744.356  | 8805.287  | 8677.574  | 3143.8   | 8745.821  |
| Std.     | 3332.744  | 3362.163  | 3298.311  | 1551.37  | 3332.058  |
| min      | 4218.81   | 4293.22   | 4209.76   | 706.88   | 4266.84   |
| 25%      | 5776.33   | 5821.95   | 5769.39   | 1908.94  | 5793.83   |
| 50%      | 7829.03   | 7867.15   | 7791.98   | 2318.76  | 7833.27   |
| 75%      | 11573.14  | 11699.63  | 11476.66  | 4416.84  | 11590.78  |
| max      | 16120.92  | 16212.23  | 16017.23  | 11621.19 | 16057.44  |
| variance | 11107186  | 11304139  | 10878852  | 2406747  | 11102609  |



Fig. 4.   Data set division into the train and test.

### B. Comparative Analysis

The algorithms BBO, GWO, and SMA, which were chosen for comparison in this research, are widely acknowledged for their efficacy and efficiency in a range of optimization tasks, including financial market analysis. The selection of these algorithms is based on their established history of success in handling intricate optimization tasks, which is consistent with the aims of our research. Relevant benchmarks, these algorithms have been extensively implemented and cited in prior research for comparable applications. Their widespread usage and recognition within the scientific community signifies their acceptability and validation, thereby furnishing a strong foundation for comparative analysis. The exclusion of other algorithms is warranted due to the fact that our research is primarily concerned with algorithms that have demonstrated exceptional potential in the analysis of financial markets. Furthermore, an excessive number of algorithms may compromise the precision and comprehensibility of the comparative analysis.

In order to accomplish the crucial objective of forecasting the Nasdaq index, an identical dataset was used by each model. In this article, a meticulous analysis and evaluation of the outcomes of each model were conducted to provide a complete and informative comparison of their respective performances. Establishing a comprehensive and equitable comparison requires the clarification of performance indicators used for evaluating the models. By employing a diverse set of significant metrics, as explained in the methodology section, the models were assessed. A comprehensive evaluation of the performance of each model can be conducted by using a variety of indicators, thus facilitating the determination of the model that most effectively meets the established requirements. All the many nuances of how each model performed are displayed in a thorough Table III with the findings.

TABLE III.    THE OUTCOMES OF THE  MODEL'S EVALUATION CRITERIA AND TIME COMPUTING

| Train/Test | Metrics | SVR | BBO-SVR | GWO-SVR | SMA-SVR |
|---|---|---|---|---|---|
| **TRAIN SET** | $R^2$ | 0.980 | 0.985 | 0.989 | 0.992 |
| | RMSE | 417.512 | 356.900 | 303.741 | 264.583 |
| | MAPE | 5.482 | 4.240 | 3.028 | 1.763 |
| | MAE | 382.054 | 329.413 | 225.660 | 158.062 |
| | RSE | 590.803 | 505.028 | 429.810 | 374.396 |
| | MSE | 174316.292 | 127377.539 | 92258.490 | 70004.188 |
| **TEST SET** | $R^2$ | 0.972 | 0.981 | 0.988 | 0.991 |
| | RMSE | 265.540 | 217.667 | 173.450 | 149.248 |
| | MAPE | 1.666 | 1.376 | 1.073 | 0.930 |
| | MAE | 215.186 | 174.385 | 134.495 | 116.260 |
| | RSE | 376.411 | 308.550 | 245.870 | 211.564 |
| | MSE | 70511.446 | 47379.041 | 30084.774 | 22275.026 |
| **Time** | Second | 0.156 | 163.85 | 221.81 | 139.94 |

Initially, the obtained outcome was used in the selection process of the SVR model. The decision to develop the SVR model was made after a comprehensive review of the data due to its exceptional performance. Between the commencement of 2015 and the midpoint of 2023, the Nasdaq index data underwent a procedure that included the selection of relevant data and normalization. The collection of important insights will benefit the decision-making process using this detailed technique. The evaluation score for SVR alone is 0.972 in $R^2$, which, as shown in Table III, has increased due to improvements in the problematic optimizers. The $R^2$ criteria values for the BBO, GWO, and SMA algorithms are 0.981, 0.988, and 0.991, respectively, indicating the possibility of selecting the optimal course of action. When compared to other optimizers, superior results are produced by the SMA optimizer. The findings of the RMSE model shown in Table III also confirm the superiority of the SMA optimizer. The RMSE values for SVR, BBO-SVR, GWO-SVR, and SMA-SVR are 265.54, 217.667, 17.3450, and 149.248, respectively. Furthermore, the hybrid models outperform the SVR model, suggesting that tweaking the model's hyperparameters can be beneficial for optimization. However compared to SVR, hybrid models require longer running times. The rationale is that while hyperparameters of the SVR are manually set, hybrid-model hyperparameters are optimized via metaheuristic algorithms. A consistent upward trend in all metrics is observed in the training set, progressing from the SVR model to the SMA-SVR model. More precisely, the $R^2$ value, which indicates the extent to which the independent variables account for the variability of the dependent variable's variance, increases from 0.980 in the SVR model to 0.992 in the SMA-SVR model. This suggests that the SMA-SVR model exhibits superior predictive capability regarding the outcomes. The RMSE, an indicator of prediction error standard deviation, exhibits a substantial reduction from 417.512 in the SVR model to 264.583 in the SMA-SVR model. This pattern is similarly evident in MAPE and MAE, where SMA-SVR exhibits a significant decline in errors, signifying its exceptional precision in prognostication. In addition, the reduction in RSE and MSE provides additional evidence that, among the four models evaluated for the train set, SMA-SVR exhibits the most accurate predictions and the lowest error rates. Similar results are observed when the SMA-SVR model is applied to the test set. The model attains the maximum $R^2$ value of 0.991, indicating an exceptional capacity for variance prediction. The SMA-SVR model exhibits the lowest values of RMSE, MAPE, and MAE, which further validate its exceptional accuracy and dependability when forecasting on unseen data. The comparatively reduced RSE and MSE values of SMA-SVR, in contrast to the alternative models, provide additional validation for its status as the most precise and dependable model across training and testing scenarios. Fig. 5 shows evaluation of the suggested model's performance in comparison to other models during training.



Fig. 5.  Evaluation of the suggested model's performance in comparison to other models during training.

Fig. 6. Evaluation of the suggested model's performance in comparison to other models during testing.

The outcomes of the conducted tests are illustrated in Fig. 7 and Fig. 8, showcasing a robust connection between the model and the empirical data. Among the models tested, the SMA-SVR model exhibited superior performance in comparison to the individual SVR, BBO-SVR, and GWO-SVR models. It is worth mentioning that the utilization of the optimizer approach resulted in a substantial enhancement in the performance of the SVR model. Fig. 6 and Fig. 7 provide a comprehensive examination of the four models, therefore validating the superior performance of the chosen model. These results indicate that the SMA-SVR model is a potential method for precisely forecasting the intended results in the context. These results indicate that the SMA-SVR model is a potential method for precisely forecasting the intended results in the context. The proposed model has the following limitations:

The high R-squared value of 0.991 may potentially signify the occurrence of overfitting in the model. This is the result of an excessive learning load imposed by the training data, which may include noise and outliers, and may consequently hinder the model's performance when applied to novel, unseen data.

# TRAIN



Fig. 7. Evaluation metrics outcomes for the attending models during train.

**TEST**



Fig. 8. Evaluation metrics outcomes for the attending models during the test.

The impact of external variables on market volatility: The model may not comprehensively capture the stock market's vulnerability to various elements, including economic indicators, political occurrences, and market sentiment. This constraint serves to underscore the intrinsic indeterminacy of the market.

The effectiveness of the model is widely recognized to be significantly contingent upon the provision of high-quality and exhaustive data. The predictive accuracy may be substantially affected by constraints in the data, including biases or incomplete historical information.

The augmentation of computational demands and complexity is an acknowledged consequence of integrating support vector regression with the slime mould algorithm. Potentially reducing the efficacy of real-time predictions, this could require additional processing time and power.

Market-Specific Variability in Generalizability: It has been noted that the efficacy of the model may differ when applied to various financial instruments or stock markets. This feature highlights the necessity for additional verification and assessment.

It is observed that the performance of the model is highly susceptible to changes in the hyperparameters it is configured with. Potentially inconsistent performance across a variety of datasets may result from the intricacy associated with fine-tuning these parameters.

Restriction on the Study's Scope: The study's concentration on particular algorithms and optimization techniques may have overlooked alternative approaches that have yet to be investigated but could have proven to be more effective.

Concerns Regarding Regulatory Compliance, Transparency, Fairness, and Ethics: The application of sophisticated predictive models in the realm of stock trading gives rise to ethical and regulatory issues. Careful attention should be given to these critical considerations.

Further research could be focused on a number of critical domains in light of the study's findings and methodologies:

Algorithm Improvement: To further enhance the accuracy and adaptability of predictions, further improvements could be made to the integration of support vector regression and the slime mould algorithm through the investigation of additional parameters or alternative configurations.

Conducting comparative analyses with novel or less prevalent optimization algorithms could provide fresh perspectives on fluctuating market conditions by evaluating the performance of the present model.

Analyses of Real-Time Data: The model's practical efficacy and resilience in authentic market situations could be evaluated through its implementation in a real-time trading environment.

Additional Market Segments: To assess the model's adaptability across various market sectors, it might be advantageous to expand its scope to encompass a more extensive array of financial instruments, including bonds, commodities, and cryptocurrencies.

The model's performance in diverse economic and regulatory environments could be better comprehended through cross-market analysis, which involves the application of the model to stock markets in different countries or regions.

Enhancement through Inclusion of Supplementary Data Sources: The predictive accuracy of the model could potentially be improved through the integration of alternative data sources, such as social media trends, news sentiment analysis, or economic indicators.

The potential for enhanced prediction capabilities and the ability to process more intricate data patterns could be investigated through the incorporation of deep learning techniques with the existing model.

The model could be rendered accessible to a wider spectrum of users, including individuals lacking technical expertise, through the development of a user interface that is intuitive and easy to use.

In order to aid investors in comprehending and alleviating potential losses, risk management functionalities might be integrated into the model.

An evaluation of the model's consistency and dependability over prolonged durations could be accomplished through the implementation of a long-term performance analysis.

The examination of the model's reaction to market anomalies or exceptional circumstances, such as unanticipated global events or financial crises, may have valuable implications.

An investigation into the feasibility of attaining a more precise or resilient prediction outcome might involve the implementation of ensemble techniques, wherein the current model is combined with additional predictive models.

Further investigation could be warranted into the ethical and regulatory ramifications associated with the utilization of sophisticated AI for forecasting the stock market, with a specific focus on the principles of trading transparency and fairness.

## IV. CONCLUSION

By employing stock prediction techniques to evaluate asset values and ascertain prevailing market trends, a substantial competitive advantage can be gained by both individual and institutional investors. Informed decisions regarding purchasing, selling, or retaining stocks can be made by investors through the utilization of historical data and advanced algorithms. The implementation of this particular method holds significant importance for investors committed to making prudent investment decisions, as it effectively mitigates risks and increases the likelihood of achieving profitable outcomes. Many predictive algorithms and data sources were employed in this research to evaluate the complex and ever-changing domain of stock prediction. The findings suggest that the accuracy of forecasting might potentially be enhanced through the utilization of a hybrid model or an ensemble technique. Finally, the construction and evaluation of the prediction model underscored the need to use data-driven insights to establish dependable decision-making processes. This highlights the benefits of a data-centric strategy in the modern, rapidly changing business environment, as well as the possible applications of predictive analytics across a wide variety of sectors. The objective of this research was to develop models with enhanced predictive capabilities for stock prices, enabling traders and investors to use these algorithms to make informed decisions on the optimal timing and price for purchasing stocks.

In this study, the following findings were reached:

- First, the data preparation and normalization process were completed, potentially impacting how the prediction model is presented. The chosen model was then prepared to begin its data analysis.

- The best model was selected, findings were analyzed, and hyperparameters were adjusted to increase the model's effectiveness.

- The identification of the most accurate optimization technique as the primary optimizer of the model was achieved through a comparative analysis of the outcomes produced by several optimization algorithms. Among the three optimization algorithms, namely BBO, GWO, and SMA, the SMA technique exhibited the highest performance in terms of the $R^2$ evaluation criterion, with a result of 0.991, surpassing the results of 0.981 and 0.988 achieved by BBO and GWO, respectively.

### REFERENCES

[1] S. Claessens, J. Frost, G. Turner, and F. Zhu, "Fintech credit markets around the world: size, drivers and policy issues," BIS Quarterly Review September, 2018.

[2] W. Li et al., "The nexus between COVID-19 fear and stock market volatility," Economic research-Ekonomska istraživanja, vol. 35, no. 1, pp. 1765–1785, 2022.

[3] Z. Wang et al., "Measuring systemic risk contribution of global stock markets: A dynamic tail risk network approach," International Review of Financial Analysis, vol. 84, p. 102361, 2022.

[4] Z. Li, W. Cheng, Y. Chen, H. Chen, and W. Wang, "Interpretable click-through rate prediction through hierarchical attention," in Proceedings of

the 13th International Conference on Web Search and Data Mining, 2020, pp. 313–321.

[5] B. Mahesh, "Machine learning algorithms-a review," International Journal of Science and Research (IJSR).[Internet], vol. 9, no. 1, pp. 381–386, 2020.

[6] S. Ray, "A quick review of machine learning algorithms," in 2019 International conference on machine learning, big data, cloud and parallel computing (COMITCon), IEEE, 2019, pp. 35–39.

[7] E. S. Olivas, J. D. M. Guerrero, M. Martinez-Sober, J. R. Magdalena-Benedito, and L. Serrano, Handbook of research on machine learning applications and trends: Algorithms, methods, and techniques: Algorithms, methods, and techniques. IGI global, 2009.

[8] A. E. L. Bilali, A. Taleb, M. A. Bahlaoui, and Y. Brouziyne, "An integrated approach based on Gaussian noises-based data augmentation method and AdaBoost model to predict faecal coliforms in rivers with small dataset," J Hydrol (Amst), vol. 599, p. 126510, 2021.

[9] G. Sonkavde, D. S. Dharrao, A. M. Bongale, S. T. Deokate, D. Doreswamy, and S. K. Bhat, "Forecasting Stock Market Prices Using Machine Learning and Deep Learning Models: A Systematic Review, Performance Analysis and Discussion of Implications," International Journal of Financial Studies, Vol 11, Iss 94, p 94 (2023), Jan. 2023, doi: 10.3390/ijfs11030094.

[10] N. Ashfaq, Z. Nawaz, and M. Ilyas, "A comparative study of Different Machine Learning Regressors For Stock Market Prediction," 2021. doi: 10.48550/arxiv.2104.07469.

[11] S. C. Agrawal, "Deep learning based non-linear regression for Stock Prediction," IOP Conference Series: Materials Science and Engineering ; volume 1116, issue 1, page 012189 ; ISSN 1757-8981 1757-899X, 2021, doi: 10.1088/1757-899x/1116/1/012189.

[12] G. James, D. Witten, T. Hastie, R. Tibshirani, and J. Taylor, "Linear regression," in An Introduction to Statistical Learning: With Applications in Python, Springer, 2023, pp. 69–134.

[13] P. Chhajer, M. Shah, and A. Kshirsagar, "The applications of artificial neural networks, support vector machines, and long–short term memory for stock market prediction," Decision Analytics Journal, vol. 2, no. November 2021, p. 100015, 2022, doi: 10.1016/j.dajour.2021.100015.

[14] S. B. Kotsiantis, "Decision trees: a recent overview," Artif Intell Rev, vol. 39, pp. 261–283, 2013.

[15] L. Breiman, "Random forests," Mach Learn, vol. 45, pp. 5–32, 2001.

[16] R. E. Wright, "Logistic regression.," 1995.

[17] A. Natekin and A. Knoll, "Gradient boosting machines, a tutorial," Front Neurorobot, vol. 7, p. 21, 2013.

[18] J. G. De Gooijer and R. J. Hyndman, "25 years of time series forecasting," Int J Forecast, vol. 22, no. 3, pp. 443–473, 2006.

[19] W. S. Noble, "What is a support vector machine?," Nat Biotechnol, vol. 24, no. 12, pp. 1565–1567, 2006.

[20] H. Rezaei, O. Bozorg-Haddad, and X. Chu, "Grey wolf optimization (GWO) algorithm," Advanced optimization by nature-inspired algorithms, pp. 81–91, 2018.

[21] S. Mirjalili, S. M. Mirjalili, and A. Lewis, "Grey wolf optimizer," Advances in engineering software, vol. 69, pp. 46–61, 2014.

[22] D. Simon, "Biogeography-based optimization," IEEE transactions on evolutionary computation, vol. 12, no. 6, pp. 702–713, 2008.

[23] S. Saremi, S. Mirjalili, and A. Lewis, "Grasshopper optimisation algorithm: theory and application," Advances in engineering software, vol. 105, pp. 30–47, 2017.

[24] S. Li, H. Chen, M. Wang, A. A. Heidari, and S. Mirjalili, "Slime mould algorithm: A new method for stochastic optimization," Future Generation Computer Systems, vol. 111, pp. 300–323, 2020.

[25] S. Mirjalili, "Moth-flame optimization algorithm: A novel nature-inspired heuristic paradigm," Knowl Based Syst, vol. 89, pp. 228–249, 2015.

[26] A. Hadidi and A. Nazari, "Design and economic optimization of shell-and-tube heat exchangers using biogeography-based (BBO) algorithm," Appl Therm Eng, vol. 51, no. 1–2, pp. 1263–1272, 2013.

[27] A. J. Smola and B. Schölkopf, "A tutorial on support vector regression," Stat Comput, vol. 14, pp. 199–222, 2004.

[28] E. H. Houssein, M. Dirar, L. Abualigah, and W. M. Mohamed, "An efficient equilibrium optimizer with support vector regression for stock market prediction," Neural Comput Appl, vol. 34, no. 4, pp. 3165–3200, 2022, doi: 10.1007/s00521-021-06580-9.

[29] A. J. Smola and B. Schölkopf, "A tutorial on support vector regression," Stat Comput, vol. 14, pp. 199–222, 2004.

# Presenting an Optimized Hybrid Model for Stock Price Prediction

Liangchao LIU

College of Economics and Management, Jiaozuo University, Henan Jiaozuo, 454000, China

*Abstract*—In the finance sector, stock price forecasting is deemed crucial for traders and investors. In this study, a detailed comparison and analysis of various machine learning models for stock price forecasting were undertaken. Historical stock data and an array of technical indicators were utilized in these models. The enhancement of the Histogram-Based Gradient Boosting (HGBR) method for predicting the Nasdaq stock index was the focus. Optimization techniques such as genetic algorithm optimization, biologically-based optimization, and the grasshopper optimization algorithm were applied. Among these, the most promising results were shown by the grasshopper optimization method. The optimized HGBR models, namely GA-HGBR, BBO-HGBR, and GOA-HGBR, were found to have achieved significant improvements, with coefficient of determination values of 0.96, 0.98, and 0.99, respectively. These figures underscore the substantial advancement of these models as compared to the baseline HGBR model. Metrics such as Mean Absolute Error, Root Mean Square Error, Mean Absolute Percentage Error, and the Coefficient of Determination were employed to assess the performance of the models.

*Keywords*—*Stock prediction; machine learning approaches; ensemble learning; grasshopper optimization; histogram-based gradient boosting*

## I. INTRODUCTION

The task of predicting stock prices is undeniably challenging, primarily due to the inherent long-term uncertainties involved [1]. The traditional market hypothesis suggests that stock prices are unpredictable and random, but current technical analysis has revealed that previous records hold valuable information that can assist in predicting future stock values [2]. Furthermore, factors such as political developments, general economic conditions, commodity prices, investor expectations, and other stock market movements can also have a significant impact on the stock market [3]. High market capitalization is utilized to calculate stock group values, and a variety of technical factors can be used to generate statistical data on stock prices [4]. Therefore, it is essential to consider all of these factors when attempting to predict stock prices accurately. Inherent challenges arise when attempting to anticipate stock values because many traditional techniques for doing so rely on stagnant trends.

Furthermore, because there are so many variables at play, forecasting stock values is inherently difficult. The market operates like a voting machine in the near term but more like a weighing machine in the long term, indicating the possibility of forecasting longer-term market changes [5]. Machine learning (ML) is a potent technology that includes a variety of algorithms and has been shown to improve performance in particular case studies considerably. Many people think that ML has the ability to find important information and recognize patterns in datasets [6]. Ensemble models are a machine learning strategy where common algorithms are used to handle a particular problem, in contrast to standard ML approaches, and they have consistently shown higher performance when it comes to time series prediction [7][8][9].

In the field of forecasting, utilizing ensemble approaches has been found to yield more accurate results compared to single models [10]. The reason behind this is that ensembles are able to combine the predictions of multiple models, enabling them to account for potential errors and uncertainties. One of the major challenges in machine learning is overfitting, which occurs when a model performs exceedingly well on training data but fails to generalize to new data. However, ensembles are less prone to overfitting due to their reliance on multiple base models, such as bagging and boosting, which help to mitigate the risk of overfitting [11]. These techniques work by creating multiple models and combining their predictions, thereby reducing the likelihood of a single model overfitting to the training data. Ultimately, the use of ensemble approaches in forecasting can lead to more reliable and accurate predictions [12]. To conduct a comparative analysis of cutting-edge machine learning methods for forecasting stock market returns, ten years of daily historical data pertaining to the top ten equities on the Casablanca Stock Exchange were utilized. When Bilal et al. [13] used an ensemble learning approach was utilized to train six classifiers (ridge regression, LASSO regression, support-vector machine, k-nearest neighbors, random forest, and adaptive boosting) to forecast price directions one day, one week, and one month in advance. In contrast to other models, support vector machines, random forests, and adaptive boosting exhibited superior performance in short-term predictions. Ensemble learning enhanced performance metrics across all prediction horizons by a substantial margin. Sonkavde et al. [14] investigated a range of algorithms to address challenges related to stock price prediction and classification. These algorithms comprised ensemble algorithms, deep learning, supervised and unsupervised machine learning, and time series analysis.

The model presented for forecasting the Nasdaq stock market in this work is a Histogram gradient boosting regressor (HGBR). Nasdaq is one of the major stock exchanges in the United States, particularly associated with technology and internet-based businesses, and renowned for its electronic trading platform. The HGBR is a machine-learning approach that addresses regression-related issues by combining the principles of gradient boosting with histogram-based feature

splitting. It is an adaptation of the popular Gradient Boosting Machine (GBM) technique [15]. Regression and classification are the two primary subtypes of gradient boosting, a machine-learning approach for prediction. This paradigm is intended to manage complicated and substantial difficulties as opposed to simple and small ones, in contrast to previous techniques. The gradient-boosting technique known as HGBR was created expressly to overcome regression issues. This technique is renowned for its quickness and capacity to hasten decision-tree learning. By discretizing the input variables, HGBR does this by splitting extra trees into several values [16].

Providing precise forecasts is the primary goal of prediction models. To achieve this, optimizing these models can significantly improve their accuracy, especially in sectors where even a minor increase in accuracy can have a substantial impact, such as healthcare, banking, and manufacturing [17]. Different methods and models are provided to optimize the HGBR. Some of them, like Moth flame optimization [18], Biogeography-based optimization [19] and gray wolf optimization [20], are inspired by nature. The optimization methods used to optimize for the model of this paper are genetic algorithm, biogeography-based optimization and grey wolf optimization.

The genetic algorithm is a computational optimization technique that draws inspiration from natural selection and evolution. This powerful tool is widely employed to solve or estimate a wide variety of optimization and search problems, ranging from engineering and finance to biology and physics [21]. By mimicking the process of natural selection, genetic algorithms are able to efficiently navigate complex search spaces and identify optimal solutions for a wide range of problems. In essence, this approach is based on the idea that the fittest solutions are more likely to survive and reproduce, leading to a gradual improvement in the overall quality of the solution over time [22]. Overall, the genetic algorithm is a versatile and powerful tool that has revolutionized the field of optimization and has enabled researchers and practitioners to tackle some of the most challenging problems of our time [23]. Another optimization method in this paper is biogeography-based optimization is a method of optimization that takes its cues from nature. Biogeography is the study of how organisms spread and adapt through time in various habitats [19]. BBO is used to solve numerous optimization problems in a variety of disciplines, including engineering, biology, economics, and data science. Biogeography is the scientific study of the geographical distribution of living organisms. The 1960s saw the discovery and development of the fundamental mathematical equations regulating the spread of organisms [24]. The GWO algorithm [20] is an innovative solution that finds its inspiration in the social hierarchy and hunting habits of grey wolves in the wild. This nature-inspired optimization technique has gained popularity in the fields of computational intelligence and machine learning due to its effectiveness in solving complex search and optimization problems [25]. By mimicking the social behavior of grey wolves, the GWO algorithm proves to be an efficient and effective way to tackle real-world optimization challenges. Its unique approach provides a fresh perspective on the problem-solving process, allowing for a more comprehensive and dynamic method of

finding solutions [26]. This paper makes a substantial contribution to the current research on predicting stock prices by thoroughly examining and analyzing several machine-learning algorithms. The application of optimization techniques, like genetic algorithm optimization, biologically-based optimization, and the grasshopper optimization algorithm, adds a layer of depth to the inquiry. The focus on improving the Histogram-Based Gradient Boosting technique for forecasting the Nasdaq stock index is remarkable. This study underscores the pragmatic significance for investors, emphasizing the cruciality of utilizing historical data and sophisticated algorithms to guide investment choices. The proposal to utilize ensemble techniques or hybrid models is in line with the progressive nature of stock prediction, recognizing the intricate and dynamic character of the market. The recognition of the grasshopper optimization algorithm as the most efficient optimizer contradicts current beliefs and offers a nuanced viewpoint on optimization methods in predicting stock prices. To summarize, this study enhances previous research by improving and perfecting the HGBR method, demonstrating the efficacy of particular optimization techniques, and providing practical guidance for investors in the ever-changing field of stock prediction.

According to the reviewed literatures, the main research gaps and novelties of the paper can be stated as follows.

*A. Research Gaps*

The field of stock price prediction, especially for the Nasdaq stock index, has long faced difficulties because of the complex interplay between contributing factors and the stock's intrinsic unpredictability. Due to their dependence on stagnant trends and inadequate analysis of the numerous factors influencing the stock market, traditional tactics frequently fail. The intricacy of this problem is exacerbated by the tendency of many machine learning models to overfit, which causes them to perform incredibly well on training data but poorly on new, untested data. Moreover, ensemble approaches are often underutilized in current models, despite the fact that they have been demonstrated to provide improved time series prediction accuracy by combining several predictions and reducing errors and uncertainties. Furthermore, although a number of optimization strategies, including Moth flame, Biogeography-based, and Gray Wolf optimization, have been studied in the literature, there is a dearth of thorough evaluation and comparison of these strategies, especially when it comes to using them to optimize the Histogram-Based Gradient Boosting (HGBR) method for stock price forecast-making.

*B. Novelties of the Work*

Using the powerful Histogram-Based Gradient Boosting Regressor (HGBR) in conjunction with cutting-edge optimization methods like genetic algorithms, biologically-based optimization, and most notably, the grasshopper optimization algorithm, this study presents an optimized hybrid model for the prediction of Nasdaq stock prices. The creation of the GA-HGBR, BBO-HGBR, and GOA-HGBR models is the result of a thorough comparison and empirical examination of various optimization techniques, which is where the innovation lies. These models have amazing coefficients of determination values and show significant improvements over

the baseline HGBR model. The exceptional efficacy of the grasshopper optimization method is revealed in this study, which is innovative in its application to stock price prediction. Furthermore, the huge dataset that was obtained from Yahoo Finance and the Nasdaq Stock Exchange and included a wide range of factors over a long period of time offers a distinctive and reliable basis for the predictive analysis. By lowering the chance of overfitting, this study not only fills in the holes in the current predictive models but also offers fresh perspectives on how well ensemble and nature-inspired optimization strategies might improve stock price predictions.

Lastly, the paper's structure is broken down into multiple sections, each of which focuses on a different aspect of the in-depth investigation that was done:

In Section II, the research methodology is the main topic of discussion in this section. It includes the explanation of the data that was utilized, the specifics of the model that was used, the optimization strategies that were used, and the evaluation criteria. The purpose of Section III is to present the study's result and discussion. Finally, Section IV concludes the paper.

## II. METHODOLOGY

### A. Data Description

This data was acquired from the Yahoo Finance Website to compile a complete historical dataset of publicly listed firms. A broad spectrum of valuable data, encompassing daily stock prices and trading volumes, was made available by this source. Five important variables—Open, High, Low, Close prices, and Trading Volume—were the main focus of the analysis in this work in order to train and test our model. To comprehend the dynamics of the stock market, these factors are essential. The opening price of a stock or other financial instrument is the price at which it is traded at the start of the trading day. It establishes the foundation for every trading day. The stock price may fluctuate and reach its highest point, referred to as the High price, during the trading day. This represents the day's peak demand or valuation. On the other hand, the price can potentially fall to what is known as the Low price—its lowest point of the day. This represents the lowest demand or valuation. The closing price is the last trading price at the conclusion of the day. It is frequently used as a benchmark for the day's performance of the stock. Additionally, trading volume is the total number of shares, contracts, or financial instruments that are exchanged in a given trading day. Elevated volumes may suggest heightened attention or involvement in a specific stock. Additionally, this dataset was supplemented with data collected directly from the Nasdaq Stock Exchange, ensuring its comprehensiveness. Access to supplementary trade metrics and market indicators was granted by this esteemed financial data source, enhancing our understanding of market dynamics. The dataset, which spans a significant timeframe from January 2015 to June 2023, includes various market conditions, including periods of stability and volatility, owing to its wide temporal range. Given its diverse array of data points that can be harnessed to construct a comprehensive industry portrait, it can be claimed that this dataset was ideally suited for robust model training and evaluation.

Nasdaq, a preeminent global stock exchange established in 1971, Nasdaq has come to represent innovation, technology, and sophisticated financial markets. Its inception aimed to modernize and streamline the stock trading process, introducing revolutionary changes to the conventional open-outcry system [27].



Fig. 1. Data division into training and testing.

The training and testing portions of the produced dataset are separated, as illustrated in Fig. 1. Data analysis and machine learning both start with the division of data into training and test sets. You may evaluate the results of the model and generalization skills using this technique.

### B. Description of the Applied Model

*1) Histogram-based gradient boosting:* HGBR represents a subtype of Gradient Boosting Regressor that accelerates the computation of the gradients and Hessians of the loss function by using histograms [28]. The algorithm starts by fitting a regressor to the training data, and then it fits other regressors to the residual errors of the first ones [15]. Weak learners are weighted together to form the final algorithm. The algorithm's main goal is to reduce the loss function:

$$L = \sum_{i=1}^{N} (y_i - \hat{y}_i)^2 \qquad (1)$$

The approach fits a weak learner $h_t(x)$ to the residual errors of the prior regressors at each iteration. The decision tree used by the weak learner divides the data into bins according to the values of the input characteristics. The approach then directly determines, rather than estimating, the gradients and Hessians of the loss function using the histogram data. Then, using precise gradients and Hessians, the weight of the learner is determined. Understanding categorical characteristics and values that are missing organically by generating new bins for each category or missing value is one benefit of histogram gradient boosting. The ultimate model is a weighted average of each weak learner separately:

$$\hat{y}(x) = \sum_{t=1}^{T} \alpha_t h_t(x) \qquad (2)$$

where, $\alpha_t$ is the learner's weight for the $t$-th weak learner.

### C. Obtained Optimization Method

When it comes to developing a reliable machine learning model, one of the most critical factors is optimizing the hyperparameters correctly. This can significantly impact the accuracy, precision, and recall of the model, which in turn can have a significant impact on the quality of the decisions and forecasts that it produces [29]. By taking the time to fine-tune the hyperparameters, developers can ensure that their model performs optimally and is better equipped to handle real-world scenarios [30].

*1) Genetic algorithm:* GA is an algorithm used for solving optimization and search issues that replicates the process of natural selection. Its basic principle is to repeatedly apply genetic operators like selection, crossover (recombination), and mutation on a population of candidate solutions or individuals to produce new individuals. The fitness function, which gauges the caliber of the solution, is then used to assess the new people. Until a workable solution is identified, this procedure is repeated over several generations [22].

*a) GA consists of three key elements [31]:* Each person is represented by a chromosome, which is a string of numbers or letters. The exact issue being handled determines the encoding. Evaluation of the fitness function is used to gauge each person's contribution to the quality of the solution. The fitness feature was created with the current issue in mind.

Using evolutionary operators, new individuals can be produced from existing ones. Selection, crossover, and mutation are the three most often utilized operators. To choose the most fertile people, selection is utilized. Chromosomes from two people can be combined through a process called crossover to create a third person. The mutation is utilized to induce minor, random alterations in an individual's chromosomes. It's essential to remember that GA is a heuristic optimization technique; it cannot be relied upon to discover the best overall solution, but it can offer a good one at a reasonable computing cost. However, for large-scale issues, it could be computationally demanding and time-consuming, particularly if the dataset is sizable and the training procedure is drawn out [32].

*2) Biological-based optimization:* BBO, a natural-inspired optimization approach, is based on the concepts of biogeography, a scientific field that investigates how species are dispersed across time in varied ecosystems. BBO is used to handle optimization difficulties in various fields, including engineering, biology, economics, and data science. Biogeography is the study of how biological organisms are distributed geographically. The discovery and development of mathematical equations that control how organisms disperse occurred in the 1960s [24]. The concept of Biogeography-Based Optimization has caught the attention of an engineer who believes that nature can teach us valuable lessons. This algorithmic approach was developed based on the principles of biogeography, which include the birth of new species, species migration between islands, and the extinction of species. Back in 2008, Dan Simon introduced this flexible and metaheuristic strategy. It uses a mathematical framework to explain how animals move across habitats, seeking refuge from unfavorable conditions and gravitating towards more hospitable ones. The Habitat Appropriateness Index is a helpful tool for evaluating and recording the suitability of different habitats. It relies solely on the objective function of the optimization problem. One of the most esteemed evolutionary algorithms is biogeography-based optimization. This algorithm systematically enhances the best solutions by optimizing a function based on a specific quality or fitness function [33].

*3) Grasshopper optimization algorithm:* The Grasshopper Optimization Algorithm, a popular metaheuristic algorithm, draws inspiration from nature. Finding the finest solutions that produce the biggest potential outcome is the key objective, and randomization is used to prevent being caught in local optima. The method has shown to be very effective and efficient in optimization thanks to its speedy convergence and impressive exploration abilities. GOA has performed better in test problems than a variety of other approaches, proving its excellence and promise in practical applications. GOA is also adaptable, balancing exploitation and exploration to ensure the optimal result is reached. This unique characteristic makes GOA an excellent choice for research applications. The overall cycle of the GOA optimizer is shown in Fig. 2.

Fig. 2.   Comprehensive cycle of GOA.

Suggested GOA, a Swarm Intelligence algorithm. [34] proposed GOA. Each grasshopper's position in the swarm, which is patterned after the behavior of grasshoppers, which regularly form swarms, represents a potential solution. The position of the $i$th grasshopper is indicated by the following equation:

$$X_i = S_i + G_i + A_i \qquad (3)$$

where, $S_i$ is for social interaction, $G_i$ stands for gravity and $A_i$ stands for wind advection.

The following equation, with the gravity element removed and the direction of the wind considered to be toward the target, states the equation adjusted for $N$ grasshopper optimization:

$$X_i^d = c\left(\sum_{\substack{j=1 \\ j \neq i}}^{N} \frac{ub_d - lb_d}{2} s\left(\left|x_j^d - x_i^d\right|\right) \frac{x_j - x_i}{d_{ij}}\right) + \widehat{T_d} \qquad (4)$$

The symbol $d_{ij}$ represents the separation among the $i$th and $j$th grasshoppers, while the function $s$ represents the strength of the social forces, where $l$ stands for the attractiveness scale and $f$ for the level of attraction, all of which are calculated using the equations below:

$$\begin{aligned} d_{ij} &= \left|d_j - d_i\right| \\ s(r) &= f e^{\frac{-r}{T}} - e^{-r} \end{aligned} \qquad (5)$$

The coefficient $c$, which decreases the comfort zone proportionately to iterations, is found using the equation.

$$c = c_{max} - l\frac{c_{max} - c_{min}}{L} \qquad (6)$$

where, $l$ is the current iteration, $C_{max}$ denotes the maximum value, $C_{min}$ denotes the minimum value, and $L$ denotes the maximum number of iterations. How the GOA optimizer works from the beginning to the end of the process is shown in Fig. 3.

### D. Evaluation Criteria

In statistics and machine learning, evaluation metrics are the chosen quantitative measurements for assessing the efficacy of prediction models. They assist in determining a model's ability to produce precise predictions on as-yet-unobserved data. The kind of prediction problem and the specific analytic goals determine the optimum assessment metric. Performance metrics, including $RMSE$, $MAE$, $MAPE$, and $R^2$ were employed in this study's predictive measures to assess the constructed forecasting models' predictive accuracy. A collection of mathematical formulas for these measurements is provided below:

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n}} \qquad (7)$$

$$MAPE = \frac{1}{n}\sum_{i=1}^{n}\left|\frac{y_i - \hat{y}_i}{y_i}\right| \qquad (8)$$

$$MAE = \frac{\sum_{i=1}^{n}|y_i - \hat{y}_i|}{n} \qquad (9)$$

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \qquad (10)$$

Fig. 3.  Flowchart of the main optimization method.

## III. RESULT AND DISCUSSION

### A. Data Statistical Results

Table I presents the statistical results of these data points, offering insights into their distribution and variability. The table indicates the count, mean, standard deviation, minimum, 25th percentile, 75th percentile, and maximum values for each variable. The count for all variables stands uniformly at 2,137, ensuring a consistent dataset for analysis. The mean values provide an average level of each variable, with the mean Close price at 8,745.821, suggesting an overall higher closing trend in the dataset. The standard deviation, particularly high in the case of High and Low prices (3,362.163 and 3,298.311 respectively), indicates significant variability and potential volatility in the market. The minimum and maximum values highlight the range of the dataset, with a notable range in the High price (from 4,293.22 to 16,212.23). The 25th and 75th percentiles reveal the distribution's skewness, where a noticeable difference is seen in the volume, and indicating periods of both low and high trading activity.

### B. Comparative Analysis

The efficacy of the models given was evaluated using a variety of standard metrics including $MAE$, $MAPE$, $R^2$, and $RMSE$. These metrics provide a thorough analysis of the forecast accuracy of the models. The performance indicators for four models, HGBR, GA-HGBR, BBO-HGBR, and GOA-HGBR, are summarized in Table II. Utilizing historical stock price information for a Nasdaq stock market index, covering from January 2015 to June 2023, these models were created and assessed.

Based on the results shown in Table II, it is clear that the GOA-HGBR model performs better than the other models in terms of predicting accuracy. The model's ability to accurately represent the complex temporal patterns and correlations contained in stock price data is demonstrated by its impressively low values for $MAE$, $MAPE$, and $RMSE$. These findings imply that the GOA-HGBR model may be a trustworthy resource for identifying potential market trends and making wise investment choices.

TABLE I. DATA STATISTICAL RESULTS

| count | 2137 | 2137 | 2137 | 2137 | 2137 |
|---|---|---|---|---|---|
| **mean** | 8744.356 | 8805.287 | 8677.574 | 3143.8 | 8745.821 |
| **Std.** | 3332.744 | 3362.163 | 3298.311 | 1551.37 | 3332.058 |
| **Min** | 4218.81 | 4293.22 | 4209.76 | 706.88 | 4266.84 |
| **25%** | 5776.33 | 5821.95 | 5769.39 | 1908.94 | 5793.83 |
| **75%** | 11573.14 | 11699.63 | 11476.66 | 4416.84 | 11590.78 |
| **max** | 16120.92 | 16212.23 | 16017.23 | 11621.19 | 16057.44 |

**Train**



Fig. 4. The results of the evaluation criteria of the developed models during training.

## TEST



Fig. 5.   The results of the evaluation criteria of the developed models during testing.

TABLE II.        THE RESULTS OF MODEL PERFORMANCE CRITERIA FOR THE NASDAQ INDEX

| MODEL / Metrics | TRAIN SET | | | | TEST SET | | | |
|---|---|---|---|---|---|---|---|---|
| | *RMSE* | *MAPE* | *MAE* | $R^2$ | *RMSE* | *MAPE* | *MAE* | $R^2$ |
| HGBR | 485.22 | 3.23 | 305.86 | 0.9726 | 337.05 | 2.18 | 274.83 | 0.9543 |
| GA-HGBR | 394.17 | 3.06 | 259.86 | 0.9819 | 275.18 | 1.83 | 221.78 | 0.9609 |
| BBO-HGBR | 324.99 | 2.86 | 231.27 | 0.9877 | 208.27 | 1.31 | 162.18 | 0.9825 |
| GOA-HGBR | 278.33 | 2.38 | 197.58 | 0.9910 | 150.97 | 0.94 | 117.44 | 0.9908 |

When the performance of the four models in Table II is compared, it can be observed that the GOA technique, followed by BBO and GA, produced the best results for optimizing the hyperparameters of the model that is being presented. The baseline HGBR model, although demonstrating robust prediction ability, acts as the standard for assessing the effectiveness of hybridization. The test set demonstrates an RMSE of 337.05 and an $R^2$ of 0.9543, providing a solid basis for evaluating the hybrid models. The incorporation of the genetic algorithm in the hybrid model results in significant enhancements. The GA-HGBR model demonstrates a decrease in the RMSE to 275.18, the MAPE to 1.83, and the MAE to 221.78 in the test set. The increase in $R^2$ (0.9609) indicates a higher level of accuracy in fitting the data, implying that the genetic algorithm successfully optimizes the hyperparameters to improve prediction accuracy. The integration of BBO into the hybrid model showcases additional enhancement. Significantly, BBO-HGBR demonstrates superior performance

in the test set, as evidenced by its lower RMSE (208.27), MAPE (1.31), and MAE (162.18), indicating an enhanced capacity to accurately capture stock price patterns. The $R^2$ value of 0.9825 confirms the effectiveness of BBO in optimizing the model for enhanced accuracy in forecasting. The outcomes of the GA-HGBR, BBO-HGBR, and GOA-HGBR findings as 0.96, 0.98, and 0.99, respectively, demonstrate the improvement in the model outcomes. The evaluation results of the developed models are shown in Fig. 4 and Fig. 5, and as it is evident in the figures, it can be seen that GOA-HGRB has the best results for all evaluation criteria. The outcomes demonstrate that the prediction result has been enhanced by the optimized model. For the $R^2$ evaluation criterion, the GOA-HGBR model, which is optimized using the GOA technique, has a result of 0.99. This result demonstrates that optimization has a beneficial impact on predicting when compared to the HGBR model, which is not optimized. Without using the optimization approach, the HGBR was 0.95.

The comparison of the developed models is illustrated in Fig. 6 and Fig. 7. The principal objective of this research endeavor was to assess the efficacy of the GOA-HGBR model in predicting NIKKEI 225 closing prices between 2013 and 2022. The comprehensive findings of this investigation are detailed in Table III, which offers an abundance of information regarding the accuracy and effectiveness of the model across various indices. With a $R^2$ value of 0.9870 for the NIIKEI 225 data set, the GOA-HGBR model exhibited superior performance compared to other models that were comparable

to the NASDQE index data set. The results suggest that the GOA-HGBR model could potentially be a valuable instrument for forecasting the forthcoming values of the aforementioned indices. Financial analysts and investors may utilize this information to aid in the formation of well-informed investment decisions. In its entirety, this research offers substantial contributions to the existing body of knowledge regarding stock price forecasting and underscores the potential of the GOA-HGBR model in predicting forthcoming financial market trends.



Fig. 6. Fit diagram of GOA-HGBR with other developed models during training.



Fig. 7. Fit diagram of GOA-HGBR with other developed models during testing.

TABLE III. THE RESULTS OF MODEL PERFORMANCE CRITERIA FOR THE NIKKEI 225 INDEX

| MODEL / Metrics | TRAIN SET | | | | TEST SET | | | |
|---|---|---|---|---|---|---|---|---|
| | RMSE | MAPE | MAE | $R^2$ | RMSE | MAPE | MAE | $R^2$ |
| HGBR | 359.5165 | 2.7144 | 288.0514 | 0.9783 | 185.5216 | 1.4299 | 145.5727 | 0.9688 |
| GA-HGBR | 335.9287 | 2.4682 | 260.1081 | 0.9821 | 175.7187 | 1.4188 | 135.3647 | 0.9712 |
| BBO-HGBR | 270.3312 | 2.1711 | 206.9693 | 0.9842 | 143.3115 | 1.4113 | 121.7996 | 0.9737 |
| GOA-HGBR | 247.8360 | 2.0554 | 186.8910 | 0.9873 | 121.9168 | 1.3252 | 98.1423 | 0.9870 |

In conclusion, the study and its findings warrant the following observations regarding future research and limitations:

- One of the model's limitations is its significant dependence on historical stock data. Particularly in the volatile and unpredictable stock market, past performance is not always indicative of future results. Consequently, this may present a constraint.

- The computational demands of sophisticated algorithms such as GA-HGBR, BBO-HGBR, and GOA-HGBR may restrict their practicality in real-time trading situations that require prompt decision-making.

- Without substantial recalibration and testing, these models may not generalize well to other stock indices or markets, despite their impressive performance for the Nasdaq stock index.

- Genetic algorithm, biologically-based optimization, and grasshopper optimization algorithm comprise the bulk of the study's attention. Alternative optimization techniques might potentially produce outcomes that are superior in quality or efficiency.

Further investigations may be warranted to examine the extent to which these models can be applied to diverse financial instruments and stock markets, thereby evaluating their adaptability and resilience.

Future Insights:

- A substantial progression would be the development of a framework for real-time data analysis and prediction, which would enable investors and traders to formulate decisions in accordance with the most up-to-date market conditions.

- Further examination of hybrid models, which amalgamate the merits of distinct algorithms, may result in the development of forecasting tools that are more precise and dependable.

- Incorporating deep learning methodologies, which have demonstrated potential in numerous predictive modeling contexts, into stock price forecasting experiments may yield novel insights and enhancements.

- Subsequent research endeavors may center on enhancing the models' capacity to navigate the frequent abrupt occurrences and market volatility that characterize the financial industry.

- By integrating these sophisticated models into user-friendly applications or platforms, they could be rendered more accessible to a wider spectrum of investors and speculators.

## IV. CONCLUSION

The best course of action for investors to take, whether to buy, sell, or hold onto stocks, can be determined by using historical data and advanced algorithms. This approach is essential for investors who are committed to making intelligent investment decisions since it lowers risks and increases the likelihood of achieving profitable results. The complex and dynamic world of stock prediction was examined in this study using a variety of predictive algorithms and data sources. These findings suggest that an ensemble technique or a hybrid model may be able to anticipate more correctly. Last but not least, the creation and evaluation of the prediction model illustrated the need for data-driven insights in order to provide trustworthy conclusions. This shows the benefits of a data-centric approach in the modern, quickly changing business environment, as well as the possible applications of predictive analytics across a wide variety of sectors. In order for interested traders and investors to utilize these algorithms to buy on the correct day and at the appropriate price, this study set out to create models that could more accurately predict stock prices.

- The study's findings both support and question previous research. Utilizing a range of metrics, including Mean Absolute Error, Root Mean Square Error, Mean Absolute Percentage Error, and the Coefficient of Determination, allows for a thorough evaluation of the model's performance. The optimized hybrid genetic algorithm-based regression models, specifically GA-HGBR, BBO-HGBR, and GOA-HGBR, demonstrate substantial enhancements, achieving a coefficient of determination value of 0.9908. This not only confirms the significance of machine learning models in predicting stock prices but also undermines conventional approaches by showcasing their superior prediction powers.

- Deciding on the best model, examining the outcomes, and then modifying its hyperparameters to enhance the performance of the previously provided model.

- To further validate the efficacy of the GOA-HGBR, these algorithms were applied to and contrasted with the NIKKEI 225 index data sets.

- By contrasting the outcomes of several optimizers, the most effective optimization has been determined as the main optimizer of the model. The GOA technique yields the best results when compared to GA, BBO, and GOA, whose $R^2$ assessment criterion scores are 0.96, 0.98, and 0.99, respectively.

REFERENCES

[1] S. Asadi, E. Hadavandi, F. Mehmanpazir, and M. M. Nakhostin, "Hybridization of evolutionary Levenberg–Marquardt neural networks and data pre-processing for stock market prediction," Knowl Based Syst, vol. 35, pp. 245–258, 2012, doi: https://doi.org/10.1016/j.knosys.2012.05.003.

[2] S. Akhter and M. A. Misir, "Capital markets efficiency: evidence from the emerging capital market with particular reference to Dhaka stock exchange," South Asian Journal of Management, vol. 12, no. 3, p. 35, 2005.

[3] K. Miao, F. Chen, and Z. G. Zhao, "Stock price forecast based on bacterial colony RBF neural network," Journal of Qingdao University (Natural Science Edition), vol. 2, no. 11, 2007.

[4] J. Lehoczky and M. Schervish, "Overview and History of Statistics for Equity Markets," Annu Rev Stat Appl, vol. 5, pp. 265–288, 2018, doi: 10.1146/annurev-statistics-031017-100518.

[5] D. Shah, H. Isah, and F. Zulkernine, "Stock market analysis: A review and taxonomy of prediction techniques," International Journal of Financial Studies, vol. 7, no. 2, 2019, doi: 10.3390/ijfs7020026.

[6] E. S. Olivas, J. D. M. Guerrero, M. Martinez-Sober, J. R. Magdalena-Benedito, and L. Serrano, Handbook of research on machine learning applications and trends: Algorithms, methods, and techniques: Algorithms, methods, and techniques. IGI global, 2009.

[7] M. Ballings, D. Van den Poel, N. Hespeels, and R. Gryp, "Evaluating multiple classifiers for stock price direction prediction," Expert Syst Appl, vol. 42, no. 20, pp. 7046–7056, 2015, doi: https://doi.org/10.1016/j.eswa.2015.05.013.

[8] M. M. Aldin, H. D. Dehnavi, and S. Entezari, "Evaluating the employment of technical indicators in predicting stock price index variations using artificial neural networks (case study: Tehran Stock Exchange)," International Journal of Business and Management, vol. 7, no. 15, p. 25, 2012.

[9] C.-F. Tsai, Y.-C. Lin, D. C. Yen, and Y.-M. Chen, "Predicting stock returns by classifier ensembles," Appl Soft Comput, vol. 11, no. 2, pp. 2452–2459, 2011, doi: https://doi.org/10.1016/j.asoc.2010.10.001.

[10] M. Zounemat-Kermani, O. Batelaan, M. Fadaee, and R. Hinkelmann, "Ensemble machine learning paradigms in hydrology: A review," J Hydrol (Amst), vol. 598, p. 126266, 2021.

[11] O. Sagi and L. Rokach, "Ensemble learning: A survey," Wiley Interdiscip Rev Data Min Knowl Discov, vol. 8, no. 4, p. e1249, 2018.

[12] S. Ardabili, A. Mosavi, and A. R. Várkonyi-Kóczy, "Advances in machine learning modeling reviewing hybrid and ensemble methods," in International conference on global research and education, Springer, 2019, pp. 215–227.

[13] A. E. L. Bilali, A. Taleb, M. A. Bahlaoui, and Y. Brouziyne, "An integrated approach based on Gaussian noises-based data augmentation method and AdaBoost model to predict faecal coliforms in rivers with small dataset," J Hydrol (Amst), vol. 599, p. 126510, 2021.

[14] G. Sonkavde, D. S. Dharrao, A. M. Bongale, S. T. Deokate, D. Doreswamy, and S. K. Bhat, "Forecasting Stock Market Prices Using Machine Learning and Deep Learning Models: A Systematic Review, Performance Analysis and Discussion of Implications," International Journal of Financial Studies, Vol 11, Iss 94, p 94 (2023), Jan. 2023, doi: 10.3390/ijfs11030094.

[15] A. Guryanov, "Histogram-based algorithm for building gradient boosting ensembles of piecewise linear decision trees," in Analysis of Images, Social Networks and Texts: 8th International Conference, AIST 2019, Kazan, Russia, July 17–19, 2019, Revised Selected Papers 8, Springer, 2019, pp. 39–50.

[16] G. Ke et al., "Lightgbm: A highly efficient gradient boosting decision tree," Adv Neural Inf Process Syst, vol. 30, 2017.

[17] S. Sun, Z. Cao, H. Zhu, and J. Zhao, "A survey of optimization methods from a machine learning perspective," IEEE Trans Cybern, vol. 50, no. 8, pp. 3668–3681, 2019.

[18] S. Mirjalili, "Moth-flame optimization algorithm: A novel nature-inspired heuristic paradigm," Knowl Based Syst, vol. 89, pp. 228–249, 2015, doi: https://doi.org/10.1016/j.knosys.2015.07.006.

[19] D. Simon, "Biogeography-based optimization," IEEE Transactions on Evolutionary Computation, vol. 12, no. 6, pp. 702–713, 2008, doi: 10.1109/TEVC.2008.919004.

[20] S. Mirjalili, S. M. Mirjalili, and A. Lewis, "Grey Wolf Optimizer," Advances in Engineering Software, vol. 69, pp. 46–61, 2014, doi: https://doi.org/10.1016/j.advengsoft.2013.12.007.

[21] B. Mohan and J. Badra, "A novel automated SuperLearner using a genetic algorithm-based hyperparameter optimization," Advances in Engineering Software, vol. 175, no. September 2022, p. 103358, 2023, doi: 10.1016/j.advengsoft.2022.103358.

[22] S. Mirjalili, "Genetic Algorithm," in Evolutionary Algorithms and Neural Networks: Theory and Applications, Cham: Springer International Publishing, 2019, pp. 43–55. doi: 10.1007/978-3-319-93025-1_4.

[23] A. Lambora, K. Gupta, and K. Chopra, "Genetic algorithm-A literature review," in 2019 international conference on machine learning, big data, cloud and parallel computing (COMITCon), IEEE, 2019, pp. 380–384.

[24] V. Garg, K. Deep, K. A. Alnowibet, H. M. Zawbaa, and A. W. Mohamed, "Biogeography Based optimization with Salp Swarm optimizer inspired operator for solving non-linear continuous optimization problems," Alexandria Engineering Journal, vol. 73, pp. 321–341, 2023, doi: https://doi.org/10.1016/j.aej.2023.04.054.

[25] H. Faris, I. Aljarah, M. A. Al-Betar, and S. Mirjalili, "Grey wolf optimizer: a review of recent variants and applications," Neural Comput Appl, vol. 30, pp. 413–435, 2018.

[26] H. Rezaei, O. Bozorg-Haddad, and X. Chu, "Grey wolf optimization (GWO) algorithm," Advanced optimization by nature-inspired algorithms, pp. 81–91, 2018.

[27] A. Abraham, B. Nath, and P. K. Mahanti, "Hybrid intelligent systems for stock market analysis," in Computational Science-ICCS 2001: International Conference San Francisco, CA, USA, May 28—30, 2001 Proceedings, Part II 1, Springer, 2001, pp. 337–345.

[28] S. Md. M. Hossain and K. Deb, "Plant Leaf Disease Recognition Using Histogram Based Gradient Boosting Classifier," in Intelligent Computing and Optimization, P. Vasant, I. Zelinka, and G.-W. Weber, Eds., Cham: Springer International Publishing, 2021, pp. 530–545.

[29] L. Yang and A. Shami, "On hyperparameter optimization of machine learning algorithms: Theory and practice," Neurocomputing, vol. 415, pp. 295–316, 2020.

[30] B. Bischl et al., "Hyperparameter optimization: Foundations, algorithms, best practices, and open challenges," Wiley Interdiscip Rev Data Min Knowl Discov, vol. 13, no. 2, p. e1484, 2023.

[31] E. Alkafaween, A. B. A. Hassanat, and S. Tarawneh, "Improving initial population for genetic algorithm using the multi linear regression based technique (MLRBT)," Communications-Scientific letters of the University of Zilina, vol. 23, no. 1, pp. E1–E10, 2021.

[32] D. M. Rocke and Z. Michalewicz, "Genetic algorithms+ data structures= evolution programs," J Am Stat Assoc, vol. 95, no. 449, p. 347, 2000.

[33] K. Cho et al., "Learning phrase representations using RNN encoder-decoder for statistical machine translation," arXiv preprint arXiv:1406.1078, 2014.

[34] S. Saremi, S. Mirjalili, and A. Lewis, "Grasshopper Optimisation Algorithm: Theory and application," Advances in Engineering Software, vol. 105, pp. 30–47, 2017, doi: https://doi.org/10.1016/j.advengsoft.2017.01.004.

# Scalable Accelerated Intelligent Charging Strategy Recommendation for Electric Vehicles Based on Deep Q-Networks

Xianhao Shen[1], Zhen Wu[2], Yexin Zhang[3], Shaohua Niu[4]*

College of Information Science and Engineering, Guilin University of Technology, Guilin, China[1, 2, 3]
Guangxi Laboratory of Embedded Technology and Intelligent System, Guilin University of Technology, Guilin, China[1, 2, 3]
School of Mechanical and Electrical Engineering, Beijing Institute of Technology, Beijing, China[4]

*Abstract*—With the rapid development of electric vehicles, their charging strategies significantly impact the overall power grid. Solving the spatiotemporal scheduling problem of vehicle charging has become a hot research topic. This paper focuses on recommending suitable charging stations for electric vehicles and proposes a scalable accelerated intelligent charging strategy recommendation algorithm based on Deep Q-Networks (DQN). The strategy recommendation problem is formulated as a Markov decision process, where the continuous sequence of regional charging requests within a time slice is fed into the DQN network as the input state, enabling optimal charging strategy recommendations for each electric vehicle. The algorithm aims to maintain regional load balance while minimizing user waiting time. To enhance the algorithm's applicability, a scalable, accelerated charging strategy framework is further proposed, which incorporates information filtering and shared experience pool mechanisms to adapt to different expansion scenarios and expedite strategy iterations in new scenarios. Simulation results demonstrate that the proposed DQN-based strategy recommendation algorithm outperforms the shortest path-first strategy, and the scalable, accelerated charging strategy framework achieves a 64.3% improvement in iteration speed in new scenarios, which helps to reduce the cloud server load and saves overheads.

*Keywords*—*Scalable acceleration; smart charging; Deep Q-network; Markov decision*

## I. INTRODUCTION

In recent years, the global energy structure has slowly transitioned towards low-carbon resources, with low-carbon energy gradually gaining a higher share in the power sector. China has also announced its efforts to achieve carbon neutrality by 2060, which will stimulate the development of the renewable energy industry in the country. According to the Renewable Energy Market Report 2023 published by the International Energy Agency, the global installed capacity of renewable energy saw an increase of over 50% in 2023 compared to the previous year, marking the most significant annual increment since 1999. By the end of 2023, the number of new energy vehicles in China reached 20.41 million, with 6.278 million charging piles available, resulting in an electric vehicle (EV) to charging pile ratio of approximately 3.3:1. The rapid growth of electric vehicles has led to an explosive demand for charging infrastructure, presenting both new challenges and opportunities for the power grid. In addition, in the face of a vast domestic user group, the existing charging stations in cities are gradually overburdened due to their sparse and uneven distribution, resulting in a severe mismatch between the current rate of new charging piles in China and the growth rate of new EV sales, and an urgent need for construction. This brings new challenges and opportunities for the intelligent charging strategy for electric vehicles.

Existing research has mainly focused on two aspects: energy storage scheduling at charging stations [1-4] and recommendation of charging strategies for electric vehicles [5-8]. Energy storage scheduling involves storing electrical energy generated by photovoltaic power generation [9-11] and managing cross-temporal energy dispatch to allocate electricity across different charging scenarios, mitigating sustained load pressure on the power grid [12]. In reference [13], a cross-temporal scheduling model integrating photovoltaic power and energy storage systems was constructed. It stored electricity during periods of low power consumption and released it during peak periods to meet changing demands. However, such cross-temporal energy scheduling algorithms rely on accurate energy usage prediction and suffer from limited energy storage efficiency and high costs. Regarding the research on the recommendation of charging strategies for electric vehicles, reference [14] developed a data-driven framework for energy prediction and utilized dynamic programming algorithms to seek optimal charging strategies. However, data-driven approaches become increasingly ineffective as the volume of data grows. Reference [15] proposes a strategy for the localization and route planning of public charging infrastructure for logistics companies based on a two-tier scheme. A two-tier genetic algorithm is used to derive the optimal routing and charging plan, and a simulated annealing descent algorithm is used to select charging station locations. The proposed method is tested and compared with a meta-heuristic approach using a benchmark instance with charging stations. Reference [16] proposes a nonlinear integer programming model with multiple objectives, including minimizing the average daily acquisition and charging costs of the electric bus routes, minimizing the time cost of waiting for charging of the electric buses and maximizing the charging revenues of the electric buses to synergistically realize the vehicle types allowed to be charged in each time window, the daily service journeys and charging journeys allocated to each electric bus. Subsequently, an algorithm was developed to

solve the formulated optimization model by combining enumeration with branching and pricing to solve the nonlinear problem. Reference [17] explored ordered charging strategies for electric vehicles using Monte Carlo algorithms, but the probabilistic nature of Monte Carlo algorithms introduces uncertainties in accurately assessing the quality of strategies. Reference [18] proposed a decision framework for charging and repositioning agent-based Shared Autonomous Electric Vehicles (SAEVs) fleets, which adjusts charging before expected demand, spatially and temporally dispersing the demand to reduce peak loads on the power grid and minimizes anticipated costs for operators. However, this framework does not consider the temporal evolution of SAEV demand and Electric Vehicle Charging Station (EVCS) supply or the cost of electricity, as its objective function only seeks to minimize response time rather than balancing charging frequency and response time.

The recommendation of charging strategies for electric vehicles is essentially a temporal scheduling problem [19], but numerous uncertain factors complicate the problem. With the rapid development of reinforcement learning, Markov decision models are well suited for charging strategy recommendations. In reference [20], a novel Markov decision process was constructed, dividing all connected electric vehicles into groups at each time step based on their charging priorities. Reinforcement learning agents were then employed to determine the charging proportions for each group of vehicles during each time interval. However, the arrival time and battery level of each electric vehicle at the charging station must be known to allocate it to a priority group. In reference [21], a graph reinforcement learning-based representation method integrates multi-dimensional information from charging stations, traffic nodes, and grid buses into a graph using feature projections. Graph convolution of coupled system states can then be implemented to facilitate environment perception. In reference [22], a novel multi-agent mean-field hierarchical reinforcement learning (MFHRL) framework was proposed to provide proactive charging and relocation advice for electric taxi drivers, maximizing the long-term cumulative rewards of their orders. The framework employed hierarchical reinforcement learning, with the manager setting goals that inherently guide the decision-making of workers, who receive rewards for following these goals. The integration of each level in the two hierarchies with mean-field approximation was carried out to incorporate the mutual influence of agents in decision-making, enabling finer temporal resolution at short intervals. In reference [23], an incentive demand response model was proposed, analyzing user behavior through reinforcement learning and subsequently guiding users to select periods with sufficient power supply. However, this approach only addresses the temporal scheduling problem, while the spatial scheduling problem remains unresolved. In reference [24], a multi-agent spatiotemporal reinforcement learning approach was introduced, altering the charging decision of electric vehicles by simulating future competitive environments using a delayed access policy. Reference [25] employed neural networks as function approximators to model user demands, training a central agent to develop charging plans for electric vehicles. None of these spatio-temporal

scheduling strategies discusses the variability of the actual environment.

Considering that China is in a period of development of charging infrastructure construction, the number of charging stations in the region is increasing, and the expandability of the scheduling strategy in the actual operation process is particularly important. Existing recommendation studies seldom consider the stability of the electric power system and the load balance while adapting to the changing environment, resulting in the charging recommendation strategy having a high maintenance and upgrading cost, and the strategy's practicality is poor.

The main contributions of this paper are mainly as follows:

*1)* To address the spatiotemporal scheduling issues in traditional electric vehicle charging strategies, a smart charging strategy recommendation algorithm based on Deep Q Network (DQN) is proposed. In this approach, the charging requests within a time slot are treated as a continuous sequence of charging request states and fed into the DQN network to generate optimal charging strategy recommendations for each electric vehicle.

*2)* To enhance the applicability of the proposed algorithm, an expandable and accelerated regional charging strategy recommendation algorithm framework is introduced. This framework utilizes a shared experience pool strategy to store strategy experiences from different regions. When a new region is added, the framework prioritizes training using experiences stored in the shared experience pool. At the same time, new experiences are stored in the experience pool of the new region. This significantly reduces the training iteration time of the model. Additionally, leveraging the experiences in the shared experience pool allows the model to converge faster and better fit the charging patterns of new regions.

The overall structure of this paper is as follows. Section II provides an introduction to the smart charging strategy recommendation model for electric vehicles based on Deep Q Network (DQN). In Section III, an expandable and accelerated regional charging strategy recommendation algorithm network framework is proposed. Section IV discusses the simulation results of the algorithm. Finally, Section V presents the conclusions and future directions for further work.

## II. Smart Charging Strategy Recommendation Algorithm for Electric Vehicles based on DQN

### A. Basic DQN Concepts

Deep Reinforcement Learning (DRL) [26] is a combination of Deep Learning (DL)[27] and Reinforcement Learning (RL)[28], which retains the ability of RL to solve policy problems. It involves the continuous interaction between an individual agent and an unknown environment, where the agent takes relevant control actions to maximize its future rewards. In theory, the value function can compute the reward value for any state and action, using methods such as Q-learning [29]. The Q-learning approach stores the state-action pairs and their corresponding rewards in a table, and when the state transitions to an environment corresponding to a table entry, the action's

reward value is obtained through table lookup. However, when there are a large number of states and actions, the computation or query time for the value or lookup function significantly increases.

The key difference between DQN and RL lies in the use of neural networks to approximate the agent's value function. Specifically, the state, s, is used as input to the neural network, and the output is the value Q(s, a) and its corresponding action, a. The Greedy(s, a) function is then combined with Qmax(Q, a) to select the best value action while maintaining a certain level of exploration. DQN calculates the current action value in a manner similar to Q-learning, using the difference between it and the output of the value neural network as the loss value. This loss value is then passed into the loss function for iterative learning. During the iterative learning process, the insertion of memories from the experience pool facilitates mixed learning, resulting in a more efficient update of the neural network.

$$y_t =$$
$$\begin{cases} R_t & \#\#end \\ R_t + \gamma max_{a_{t-1}} Q_{t-1}(s_{t-1}, a_{t-1}) & \#\#\#\#\#\#\#\#\#\#\#\#\#\#\#\#\#\#not\ end\#\#\#\#\#\#\# \#\#\#\#\#\#\#\# \end{cases}$$
(1)

$$L_\theta = E(y_t - Q_t(s_t, a_t, \theta))^2$$
(2)

Eq. (1) and Eq. (2) represent the target value, $y_t$ and the reward value, $R_t$_trespectively, at time step $t$. The learning discount rate is denoted as $\gamma$, and $Q_{t-1}(s_{t-1}, a_{t-1})$ represents the value network value at time step $t-1$. $L_\theta$ corresponds to the value of the loss function.

### B. A DQN-based Recommendation Model for Smart Charging Strategies for Electric Vehicles

This study primarily focuses on providing charging strategy recommendations for electric vehicles (EVs) at public charging stations. As illustrated in Fig. 1, EVs with charging demands within a designated area send their charging requests to a central processor. They also transmit their specific vehicle information, including current battery level and location, to the central processor. The central processor collects all the charging requests from EVs within the same time period in the area, forming a temporal sequence of charging requests. This sequence serves as the input to the Deep Q-Network (DQN) for generating optimal charging strategy recommendations for all EVs within a time slice. Considering the timeliness of charging strategy recommendations, the study employs time slicing by dividing each minute into 60 time slices, with each time slice representing 1 second. Within a time slice, the processor composes timing input vectors from the states of all requests combined with the load conditions of the charging station and charging pile information, etc., and makes the correct strategy decision for the EV through a deep reinforcement learning model to guide the EV to complete the charging, which satisfies the need to maintain the load of the regional power grid while shortening the user's waiting time.

The recommendation of charging strategies for electric vehicles (EVs) can be viewed as a Markov decision process, which involves coordinating the interaction between EVs and the regional charging environment. The goal is to guide each EV to make informed decisions regarding charging strategies

while minimizing user waiting time and the load on the regional power grid. However, treating each individual EV as the main agent does not satisfy the continuity of the state space in the Markov decision process. Therefore, a time slicing approach is adopted, where all charging requests from EVs within a time slice in the region are sorted based on their submission time, forming a continuous state space for the regional charging requests. As shown in Fig. 2 and described by Eq. (3) and Eq. (4), when the agent $C_1$ submits request $q_1$, its current location state $Qarea$, state of charge $Qsoc$, and the location information of each charging station $Larea$ are combined to form the overall state $s_{q_1}$. Subsequently, the state transitions from $s_{q_1}$ to $s_{q_2}$ as the vehicle $C_2$ request is processed. The collection of all request states $s_{q_m}$ within the time slice $t_1$ forms the aggregate state $s_{t_n}$, which serves as the input to the DQN network for training in a single episode.

$$s_{q_m} = (Qarea_{q_m}, Qsoc_{q_m}, Larea_{q_m})$$
(3)

$$s_{t_n} = \{s_{q_1}, s_{q_2}, \dots, s_{q_m}\}$$
(4)

The DQN action space employed in this study corresponds to the selection of charging stations, where EVs continuously make decisions on charging stations within a time slice, and the action space is the same for all requests. As shown in Fig. 3, the action space corresponds to different discrete charging stations. The agent can choose from the following four actions: 1、 Charging station 1, 2、 Charging station 2, 3、 Charging station 3, and 4、 Charging station 4. These charging stations are randomly distributed, and their initial charging states are also randomized. The agent is trained in various stochastic environments to cope with challenges in real-world settings.



Fig. 1. Scenario of the use of DQN-based recommendation model for smart charging strategy for electric vehicles.



Fig. 2. Time slice model diagram.

Fig. 3.   Action space diagram.

Rewards provide direct or delayed feedback to the agent's decisions, enabling the agent to continually update its decisions to maximize the rewards. Rewards quantify higher-level objectives in multi-agent reinforcement learning. Specifically, in the context of electric vehicle charging strategies, the reward is set as a composite reward to expedite the training iteration of the intelligent agent. Upon making a decision regarding the request $q_t$, the intelligent agent receives the reward functions $r_{dis}(s_{q_m}, a_{q_m})$, $r_{wait}(s_{q_m}, a_{q_m})$, and $r_{load}(s_{q_m}, a_{q_m})$ as defined in Eq. (5) and Eq. (6), respectively:

$$\begin{cases} r_{dis}(s_{q_m}, a_{q_m}) = \begin{cases} -\frac{dis_{range}}{dis_{min}} \; dis_{range} + dis_{min} < mile_{min} \\ -dis_{range} \; dis_{range} + dis_{min} > mile_{min} \end{cases} \\ r_{wait}(s_{q_m}, a_{q_m}) = -t_{wait} \\ r_{load}(s_{q_m}, a_{q_m}) = \begin{cases} \sum_{n_{charge}}^{i} \frac{1}{n_{rest}} \\ \sum_{n_{charge}}^{i} \frac{1}{n_{wait}} \end{cases} \end{cases} \quad (5)$$

$$R_t(sli_t, a_t) = \partial r_{dis}(sli_t, a_t) + \beta r_{wait}(sli_t, a_t) + \delta r_{load}(sli_t, a_t) \quad (6)$$

The variables in the equation are defined as follows: $dis_{min}$ represents the shortest distance to the charging station, $dis_{error}$ denotes the difference between the selected station and the shortest station in terms of distance. $mile_{min}$ represents the minimum remaining mileage of the vehicle based on its current condition. $t_{wait}$ represents the charging waiting time. $n_{charge}$ corresponds to the sequence number of the charging station. $n_{rest}$ and $n_{wait}$ respectively indicate the number of available charging piles at the currently selected station and the number of vehicles queuing for charging. $s_{nc}$ represents the current charging status of the station, which collectively determines the overall load of the region. Finally, $\partial$, $\beta$, and $\delta$ are discount factors.

Observation value: To give the central processor a better grasp of the global information, an observation value is set for each intelligent body. They are set as shown in Eq. (7).

$$obs_m = \{Qarea_{q_m}, fpile_{q_m}, load_{q_m}\} \quad (7)$$

In the equation, $s_m$ represents the state set composed of the position status and battery status of all electric vehicles within the current region. $f_m$ denotes the number of available charging piles in the region, while $l_m$ represents the total load of the region.

### C. Reinforcing the Learning Process

The recommended smart charging strategy for electric vehicles based on DQN is illustrated in Fig. 4. The process

begins by initializing the experience replay buffer, neural network parameters, and the initial state denoted as s in the DRL model. Subsequently, the states of all electric vehicles within the region are collected to form a state set. The charging policy network and the charging value network are separately utilized to obtain the actual reward r, value network reward $r'$, next state $s'$, and action a. These parameters are then stored in the experience replay buffer. At irregular intervals, parameters are randomly sampled from the experience replay buffer and added to the EV state set for training. Following this, the value network reward $r'$ and the actual reward $r$ are input into the loss function to train the charging value network. The EV state $s$ is updated to $s'$ and the EV state set is updated iteratively until the current training round is completed. Finally, the next LSTM model predicts the EV state, and this process continues until the training is completed, resulting in the output of the trained charging value network model.



Fig. 4.   Diagram of recommended smart charging strategies for electric vehicles.

## III.   AN ALGORITHMIC FRAMEWORK FOR RECOMMENDING REGIONAL CHARGING STRATEGIES WITH SCALABLE ACCELERATION

### A. Framework Background

Currently, electric vehicles are undergoing an incredible and rapid development, leading to a continuous increase in charging demand. To alleviate the pressure on charging load, many regions have started constructing new charging stations. However, existing charging strategies [30-32] have not addressed their scalability. Adding a new charging area and starting the training of charging recommendation strategies from scratch undoubtedly incurs additional costs. Therefore, this paper proposes a scalable and accelerated framework for regional charging strategy recommendation algorithm.

### B. Framework Scenario Analysis

The individual charging station information within a single region is presented in Fig. 5 and Fig. 6, including the operational status, available quantity of charging piles, and specific locations of the charging stations. Initially, the information filtering layer is employed to select the information from the n closest charging stations to the charging-requesting vehicle, forming a new tuple of charging station information features with a length of n. The specific value of n will be described in detail in the experimental section. Subsequently, the new tuple of features is input into the DQN network for training, ultimately providing policy recommendations. The initial input states, decisions, rewards, and other parameters for each policy recommendation are stored in the network's own experience replay buffer and the

shared experience pool of the extended framework. During the training of policy recommendation across multiple regions, when updating the policy value network, random sampling from the shared experience pool is incorporated to achieve experience sharing. This facilitates accelerated training when new regions join the shared experience pool, effectively avoiding the issue of random recommendations due to insufficient initial experience pool capacity. Furthermore, the self-experience replay buffer is continually improved during the training process. Once its capacity is full and construction is completed, the framework utilizes its own experience replay buffer. Next, the applicability of this framework will be discussed based on three extended scenarios [33].

Scenario 1: Addition of new charging piles within the region. The purpose of this algorithm is to recommend the optimal charging station. Within the algorithm environment, there is a queue of information regarding available charging piles at the charging stations. When new charging piles are added to a charging station, it simply increases the count of available charging piles, without affecting the functionality of the algorithm.

Scenario 2: Addition of new charging stations within the region. The first layer of the proposed recommendation algorithm framework filters the information of all charging stations within the region. It retains a tuple of information features with a length of n ensuring that the input dimension of the DQN network remains consistent. This, in turn, guarantees consistency in the action space dimension of the DQN network. Specifically, when a vehicle makes a request, the DQN network takes a filtered queue of n nearest charging station information as input and ultimately provides policy recommendations among these n charging stations for the vehicle.

Scenario 3: Addition of a new charging area. The shared experience pool within the proposed recommendation algorithm framework is designed to address this scenario. The new region can directly utilize the shared experience pool to accelerate training, continuously accumulate and improve its own experience replay buffer, and eventually develop its specific charging strategy.



Fig. 5. Framework of the scalable regional charging policy recommendation algorithm.



Fig. 6. Flow chart of the shared experience pool.

## IV. ALGORITHM ANALYSIS

### A. DQN-based Algorithm for Recommending Smart Charging Strategies for Electric Vehicles

*1) Experiment description:* This experiment primarily simulates the decision-making behavior of electric vehicles in a region regarding public charging stations. The region is set to a size of 2000*2000 grids, and at the beginning of each experimental round, the coordinates of the charging stations within the region, as well as the positions and coordinates of the charging requests within the region, are randomly initialized. The coordinates are used to simulate real-world latitude and longitude.

In this experiment, a comparison will be made between the DQN-based intelligent charging strategy recommendation algorithm and the nearest distance-first strategy in terms of specific performance metrics such as average charging waiting time and average regional load. The experiment involves storing the location information of the region's charging stations on a server and simulating the application scenarios of the scalable regional charging strategy recommendation algorithm framework through local-server interactions.

*2) Parameter setting:* To validate the proposed algorithm, the following experiments were conducted in the simulation environment as shown in Table I. In this algorithm, the batch size of 32 was selected for each training iteration. The learning rate *lr* of the DQN network was set to 0.01, the exploration-exploitation trade-off rate $\varepsilon$ was set to 0.9, and the discount factor $\gamma$ for the policy was set to 0.9. The experience replay buffer size was set to 100,000, and the target network was updated every 100 iterations .

TABLE I. SIMULATION ENVIRONMENT CONFIGURATION TABLE

| Configuration of experiment | Specific Parameters |
|---|---|
| CPU | Intel Core i7-11700K@ 5.0GHz |
| GPU | NVIDIA GeForce GTX 1660 Super |
| Memory | 16GB |
| Operating system | Windows10 |
| Programming environment | python3.7、pytorch 1.10.2+cu102 |
| Reinforcement learning environment | gym0.10.5 |

*3) Analysis of results:* After 5000 iterations of training, as shown in Fig. 7, where the first 1000 iterations were used for the experience replay buffer population, the reward values for the DQN-based intelligent charging strategy recommendation algorithm converged to approximately -3.5. To demonstrate the performance of the proposed algorithm, we will now discuss the simulation results in detail. Fig. 8(a) presents the load situation of the charging strategy recommendation algorithm based on DQN, while Fig. 8(b) shows the load situation for the nearest distance-first strategy. It can be observed that with the increase in time steps, our proposed algorithm exhibits some fluctuations. However, it shows an overall decreasing trend, significantly different from the nearest distance-first strategy. Based on calculations, the average load per step for the DQN-based charging strategy recommendation algorithm is 1.14, while for the nearest distance-first strategy, it is 1.20, resulting in an improvement of approximately 5.0%.



Fig. 7. Iteration diagram of the DQN-based charging policy recomme-ndation algorithm model.



Fig. 8. DQN based intelligent charging strategy recommendation algorithm for electric vehicles and the nearest distance recommended load comparison diagram.

In terms of waiting time, as shown in Fig. 9, where Fig. 9(a) represents the DQN-based electric vehicle charging strategy recommendation algorithm's ability to make correct recommendations for immediate use when all charging stations in the area are initially vacant and to reasonably schedule charging plans even when all charging stations are under load in the latter part. In contrast, Fig. 9(b) depicts the nearest distance-first strategy, which fails to make optimal charging plan arrangements from the beginning. Based on calculations, the average waiting time per step for the DQN-based charging strategy recommendation algorithm is 1.75 ms, while for the

nearest distance-first strategy, it is 1.91 ms, resulting in an improvement of approximately 8.37%. It can be observed that the DQN-based charging strategy recommendation algorithm proposed in this paper not only maintains balanced area loads but also significantly reduces users' waiting time, which helps alleviate user anxiety and enhances the user experience.



Fig. 9. DQN based intelligent charging strategy recommendation algorithm for electric vehicles and the nearest distance recommended waiting time comparison diagram.

In addition to waiting time, the distance to the recommended charging station is also a criterion for measuring the algorithm's accuracy. Fig. 10(a) and 10(b) below represent the DQN-based electric vehicle charging strategy recommendation algorithm and the nearest distance recommendation algorithm, respectively. From the figures, it can be observed that the DQN-based recommendation algorithm is primarily consistent with the nearest distance recommendation algorithm. Out of 300 testing steps, the DQN-based recommendation algorithm recommended the nearest distance priority in 121 cases. In contrast, in the remaining 179 cases, it selected other charging stations to minimize total time.



Fig. 10. DQN based intelligent charging strategy recommendation algorithm for electric vehicles and the nearest distance recommendation algorithm recommended site distance comparison diagram.

In terms of enablement, the experiments tested the DQN-based scalable EV smart charging policy recommendation algorithm model size of 1.21M with an average delay of 923ms, which has a strong real-time performance and can be applied to practical scenarios.

### B. Scalable Acceleration Algorithm for Electric Vehicle Charging Strategy Recommendation

*1) Experiment description:* This experiment mainly simulates the expansion charging scenario in the region, the iteration speed of the model will be verified separately, and the expandability in different scenarios.

*2) Parameter setting:* The experimental parameters of the network part of this experiment are consistent with the DQN strategy algorithm above, i.e. the number of samples selected for one training session is 32, the learning rate lr is 0.01, the

greed rate ε is 0.9, the discount rate γ is 0.9, the experience pool size is set to 100000, and the target network following frequency is 100. The total experience pool size in the algorithm is set to 10,000,000 and the following frequency is 10,000.

*3) Analysis of results:* To cope with increasingly complex charging scenarios, the proposed scalable and accelerated electric vehicle charging strategy recommendation algorithm in this paper filters the charging station information within a single region through an information filtering layer. It forms a new charging station information feature tuple of length n, as illustrated in Fig. 11. In order to observe the impact of n on the waiting time in the recommendation algorithm, we conducted the following experiments, and the average waiting time was minimized when n was set to 9.

To validate the feasibility of the algorithm, this study conducted simulation experiments on the following scenarios based on practical application scenarios: Scenario 1: Adding new charging poles to charging stations within a region; Scenario 2: Adding new charging stations within a region; Scenario 3: Adding new charging regions; Scenario 4: Complex real-world scenarios. In Scenario 4, the number of experimental subjects in Scenario 3 was doubled. Specific parameters are shown in Table II.

Fig. 12 presents a comparison of the average waiting time for each scenario. Fig. 12(a) shows the original waiting time graph obtained from Experiment IV (A), while Fig. 12(b), Fig. 12(c), and Fig. 12(d) correspond to Scenario 1, Scenario 2, and Scenario 3, respectively. In Scenario 1, adding new charging poles within the region provides the algorithm with more choices, resulting in a significant decrease in the average waiting time to 0.59 ms. In Scenario 2, adding new charging stations slightly reduces the average waiting time to 0.96 ms. In Scenario 3, expanding the charging region results in a decreased average waiting time of 1.69 ms, representing improvements of 66.3%, 45.1%, and 3.4%, respectively. Thus, it can be concluded that the proposed scalable and accelerated electric vehicle charging scheduling recommendation algorithm remains applicable in complex scenarios, and its performance improves as the complexity of the scenarios increases.

In terms of the load aspect, as shown in Fig. 13, Fig. 13(a) represents the original load graph, Fig. 13(b) represents the load graph for Scenario 1 with an average load reduction of 0.8, Fig. 13(c) represents the load graph for Scenario 2 with an average load reduction of 0.84, and Fig. 13(d) represents the load graph for Scenario 3 with an average load reduction of 1.12. These reductions correspond to 29.8%, 26.3%, and 1.75% improvements, respectively. The average algorithm latency for each scenario is 889 ms, 893 ms, and 897 ms, representing reductions of 3.7%, 3.2%, and 2.8%, respectively. In conclusion, the proposed scalable and accelerated electric vehicle charging scheduling recommendation algorithm in this chapter reduces user waiting time while ensuring the stability of the regional load in complex scenarios. This contributes to better revenue generation for operators.



Fig. 11. Graph of results for feature tuple length n.



Fig. 12. Comparison of the average waiting time in each scenario.

TABLE II.     COMPARISON TABLE OF SCENE PARAMETERS

| Scene Name | The original scene | Scenario 1 | Scenario 2 | Scenario 3 | Scenario 4 |
|---|---|---|---|---|---|
| Size of grid | 2000*2000 | 2000*2000 | 2000*2000 | 2000*2000 | 6000*6000 |
| Number of requests per unit time | 300 | 300 | 300 | 300 | 300 |
| Number of charging stations | 20 | 20 | 40 | 20 | 40 |
| Number of charging piles per charging station | 20 | 40 | 20 | 20 | 40 |

Fig. 13. Load comparison diagram of each scenario.



Fig. 14. Comparison diagram of training iteration speed of each scene.



Fig. 15. Scenario 3 and Scenario 4 Waiting time comparison.

To validate the applicability of the proposed algorithm in real-world complex scenarios, we introduced increased complexity to the scenario parameters, and the results are shown in Fig. 14. Fig. 14(a) depicts the iteration graph of the charging strategy recommendation algorithm based on DQN, which converges after approximately 1400 iterations due to the need for experience pool storage. Fig. 14(b) represents the iteration graph for Scenario 3, while Fig. 14(c) corresponds to the iteration graph for Scenario 4. It is evident that compared to the DQN-based charging strategy recommendation algorithm, expanding the new region in training, as proposed in this study, using the shared experience pool approach eliminates the time required for storing the experience pool. Moreover, the experiences generated by the shared experience pool, compared to those randomly selected by DQN for action generation, are more practical and accelerate the fitting of model parameters, resulting in faster model iterations. In particular, the iteration speed is improved by 64.3% (500 iterations) and 67.8% (450 iterations) for Scenario 3 and Scenario 4, respectively. This significantly reduces the load on the cloud server and saves costs. Regarding waiting time, as shown in Fig. 15, Scenario 4, with the addition of more charging stations and charging piles and an expanded map area, offers users more choices, leading to a decrease in average waiting time compared to Scenario 3, reaching 0.51ms. Complex scenarios often accompany increased model execution time. However, in the simulated experiments of this algorithm in Scenario 4, the average algorithm latency remained relatively unchanged at 901ms, as mentioned earlier. The results demonstrate that as the complexity of the application scenarios increases, this algorithm can further accelerate the model iteration speed, reduce average waiting time for users, and maintain a consistent algorithm latency, showcasing its high applicability.

## V. CONCLUSION

This paper presents an intelligent electric vehicle (EV) charging strategy recommendation algorithm based on Deep Q-Network (DQN). The algorithm utilizes Markov modeling of user-requested charging events to formulate reasonable charging plans and effectively addresses the spatial scheduling issues in traditional EV charging strategies. Considering the rapid development of charging infrastructure construction in China, we propose a scalable and accelerated regional charging strategy recommendation algorithm framework. This framework not only adapts to increasingly complex and evolving charging scenarios but also maintains a consistent algorithm latency, further accelerating the iteration of the algorithm model. Experimental results show that the algorithm can improve the efficiency of charging strategy recommendation, charging waiting time, and charging demand response speed. In contrast, the expandable and accelerated charging strategy framework improves the iterative speed by 64.3% in new scenarios, which reduces the cloud server load and saves overheads. In future work, we will further refine the hardware implementation of the algorithm to realize a more efficient, precise, and practical charging strategy recommendation algorithm. This will provide superior, efficient, and convenient charging services for EVs, positively contributing to the development of innovative urban transportation.

REFERENCES

[1] M.S. Lipu Hossain, et al, "A review of controllers and optimizations based scheduling operation for battery energy storage system towards decarbonization in microgrid: Challenges and future directions," Journal of Cleaner Production, vol. 360, pp. 132188, Aug. 2022.

[2] J. Massana, L. Burgas, S. Herraiz, J. Colomer, and C. Pous, "Multi-vector energy management system including scheduling electrolyser, electric vehicle charging station and other assets in a real scenario," Journal of Cleaner Production, vol. 380, pp. 134996, Dec. 2022.

[3] G. Jianwei, G. Fangjie, Y. Yu, et al, "Configuration optimization and benefit allocation model of multi-park integrated energy systems considering electric vehicle charging station to assist services of shared energy storage power station," Journal of Cleaner Production, vol. 336, pp. 130381, Feb. 2022.

[4] L. Desheng, Z. Adama, L. Jiantang, and Y. Hongtzer, "An energy management strategy with renewable energy and energy storage system for a large electric vehicle charging station," eTransportation, vol. 6, pp. 100076, 2020.

[5] M.A. Hannan, M.S. Mollik, A.Q. Al-Shetwi, et al, "Vehicle to grid connected technologies and charging strategies: Operation, control, issues and recommendations," Journal of Cleaner Production, vol. 339, pp. 130587, Mar. 2022.

[6] Z. Jiaqing, L. Jing, and Z. Xiaohui, "Charging navigation strategy for electric vehicles considering empty-loading ratio and dynamic electricity price," Sustainable Energy, Grids and Networks, vol. 34, pp. 100987, Jun. 2023.

[7] P. Priyadarshan, K. Kazemzadeh, and B. Prateek, "Integration of charging behavior into infrastructure planning and management of electric vehicles: A systematic review and framework," Sustainable Cities and Society, vol. 88, pp. 104265, Jan. 2023.

[8] M. S. Muhammad, Z. Shengxian, M. M. Hafiz, et al, "A study of charging-dispatch strategies and vehicle-to-grid technologies for electric vehicles in distribution networks," Energy Reports, vol. 9, PP. 1777-1806, Dec. 2023.

[9] Y. Y. Kah, C. H. Hon, and K. J. Jiří, "Solar Energy-Powered Battery Electric Vehicle charging stations: Current development and future prospect review," Renewable and Sustainable Energy Reviews, vol. 169, pp. 112862, Nov. 2022.

[10] X. Fangzhou, C. Hongkun, L. Hao, and C. Lei, "Optimal planning of photovoltaic-storage fast charging station considering electric vehicle charging demand response," Energy Reports, vol. 8, pp. 399-412, Nov. 2022.

[11] Z. Sixiang, W. Yachao, J. Zhenyu, H. Tianshu, and C. Fengming, "Research on emergency distribution optimization of mobile power for electric vehicle in photovoltaic-energy storage-charging supply chain under the energy blockchain," Energy Reports, vol. 8, pp. 6815-6825, Nov. 2022.

[12] Y. Meng, Z. Lihui, Z. Zhenli, and W. Liwan, "Comprehensive benefits analysis of electric vehicle charging station integrated photovoltaic and energy storage," Journal of Cleaner Production, vol. 302, pp.126967, Jun. 2021.

[13] K. Kouka, A. Masmoudi, A. Abdelkafi, and L. Krichen, "Dynamic energy management of an electric vehicle charging station using photovoltaic power," Sustainable Energy, Grids and Networks, vol. 24, p. 100402, Dec. 2020.

[14] Z. Yi and M. Shirk, "Data-driven optimal charging decision making for connected and automated electric vehicles: A personal usage scenario," Transportation Research Part C: Emerging Technologies, vol. 86, pp. 37–58, Jan. 2018.

[15] L. Jiale, L. Zhenbo, and W. Xuefei, "Public charging station localization and route planning of electric vehicles considering the operational strategy: A bi-level optimizing approach," Sustainable Cities and Society, vol. 87, pp. 104153. Dec. 2022.

[16] J. Jinhua, B. Yiming, and W. Linhong, "Optimal electric bus fleet scheduling for a route with charging facility sharing," Transportation Research Part C: Emerging Technologies, vol. 147, pp. 104010, Feb. 2023.

[17] L. Hongyan et al., "Research on orderly charging control strategy of electric vehicles based on Monte Carlo algorithm, " Electrical Appliance and Energy Efficiency Management Technology, no. 05, pp. 59–64, 2021.

[18] M. D. Dean, K. M. Gurumurthy, F. de Souza, J. Auld, and K. M. Kockelman, "Synergies between repositioning and charging strategies for shared autonomous electric vehicle fleets," Transportation Research Part D: Transport and Environment, vol. 108, pp. 103314, Jul. 2022.

[19] W. Ning, T. Hangqi, Z. Shunbo, and L. Yuan, "Analysis of public acceptance of electric vehicle charging scheduling based on the technology acceptance model," Energy, vol. 258, pp. 124804, Nov. 2022.

[20] N. Sadeghianpourhamami, J. Deleu, and C. Develder, "Definition and Evaluation of Model-Free Coordination of Electrical Vehicle Charging With Reinforcement Learning," IEEE Transactions on Smart Grid, vol. 11, no. 1, pp. 203–214, Jan. 2020.

[21] X. Peidong, Z. Jun, G. Tianlu, et al, "Real-time fast charging station recommendation for electric vehicles in coupled power-transportation networks: A graph reinforcement learning method," International Journal of Electrical Power & Energy Systems, vol. 141, pp. 108030, Oct. 2022.

[22] E. Wang et al, "Joint Charging and Relocation Recommendation for E-Taxi Drivers via Multi-Agent Mean Field Hierarchical Reinforcement Learning," IEEE. Trans. Mob. Comput., vol. 21, no. 4, pp. 1274–1290, Apr. 2022.

[23] L. Xin et al, "Orderly charging control of electric vehicles based on MDP and incentive demand response," Journal of Electric Power Science and Technology, vol. 36, no. 05, pp. 79–86, 2021.

[24] W. Zhang et al, "Intelligent Electric Vehicle Charging Recommendation Based on Multi-Agent Reinforcement Learning," WWW, 2021.

[25] F. Tuchnitz, N. Ebell, J. Schlund, and M. Pruckner, "Development and Evaluation of a Smart Charging Strategy for an Electric Vehicle Fleet Based on Reinforcement Learning," Applied Energy, vol. 285, p. 116382, Mar. 2021.

[26] K. Arulkumaran, M.P. Deisenroth, M. Brundage, and A.A. Bharath, "Deep Reinforcement Learning: A Brief Survey," IEEE Signal Processing Magazine, vol. 34, no. 6, pp. 26–38, Sep. 2017.

[27] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," Nature, vol. 521, no. 7553, Art. no. 7553, May 2015, doi: 10.1038/nature14539.

[28] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," The International Journal of Robotics Research, vol. 32, no. 11, pp. 1238–1274, Sep. 2013.

[29] C. J. C. H. Watkins and P. Dayan, "Q-learning," Mach Learn, vol. 8, no. 3, pp. 279–292, May 1992.

[30] P. Xu, J. Zhang, T. Gao et al, "Real-time fast charging station recommendation for electric vehicles in coupled power-transportation networks: A graph reinforcement learning method," International Journal of Electrical Power & Energy Systems, vol.141, p. 2108030, Oct. 2022.

[31] X. Qiang, C. Zhong, Z. Ziqi, H. Xueliang, and L. Xiaohui, "Route Planning and Charging Navigation Strategy for Electric Vehicles Based on Real-time Traffic Information and Grid Information," IOP Conf. Ser.: Mater. Sci. Eng., vol. 752, no. 1, p. 012011, 2020.

[32] J. Zhong, N. Yang, X. Zhang, and J. Liu, "A fast-charging navigation strategy for electric vehicles considering user time utility differences," Sustainable Energy, Grids and Networks, vol. 30, p. 100646, Jun. 2022.

[33] S. Borozan, S. Giannelos, and G. Strbac, "Strategic network expansion planning with electric vehicle smart charging concepts as investment options," Advances in Applied Energy, vol. 5, p. 100077, Feb. 2022.

# Geospatial Pharmacy Navigator: A Web and Mobile Application Integrating Geographical Information System (GIS) for Medicine Accessibility

Mia Amor C. Tinam-isan[1], Sherwin D. Sandoval[2], Nathanael R. Neri[3], Nasrollah L. Gandamato[4]
Department of Information Technology-CCS, MSU-Iligan Institute of Technology, Iligan City, Philippines[1]
MSU-Iligan Institute of Technology, Iligan City, Philippines[2, 3, 4]

*Abstract*—This project introduces a web and mobile application that integrates Geographic Information Systems (GIS) to identify pharmacies with available prescription drugs, addressing the expanding role of Information and Communication Technology (ICT) in healthcare. The primary objective is to offer the general public an easy-to-use platform that locates the closest pharmacy having the searched drugs or medicines. Adopting the Rapid Application Development methodology ensures continuous engagement with stakeholders, allowing developers to closely align the application with user requirements. Essential elements of the web platform include chat functionality, inventory management, pharmacy oversight, and the display of medication listings. General users may check medication lists, search pharmacies, find pharmacy locations and the best routes, search for specific medications, access comprehensive medication information, and more with the mobile application. Fifty respondents, comprising five pharmacists and forty-five general users, expressed overall satisfaction with the system's functionality, emphasizing its ease of use and straightforward navigation across most features. This project not only amplifies the importance of ICT in the healthcare industry, but it also shows how technology can be successfully integrated to improve accessibility and expedite healthcare procedures for both the general public and professionals.

*Keywords*—*ICT in health; mobile application; web application; GIS; pharmacy mapping*

## I. INTRODUCTION

Access to necessary medical care, especially prescription medication, is an essential aspect of public health, yet it is still a major problem in the Philippines. People seeking necessary medications face a complicated web of difficulties due to the convergence of factors like escalating medicine costs, insufficient delivery methods, and limited accessibility. As stipulated in international human rights agreements, the World Health Organization (WHO) states that guaranteeing access to necessary medications is a crucial part of the right to health [1].

But in the case of the Philippines, the difficulty in obtaining prescription drugs is made worse by the lack of reliable ICT (information and communication technology) tools that can help people find prescribed medication. Prescription medicine accessibility in the Philippines is impeded by a number of variables, such as geographic differences, financial limitations, and a disjointed healthcare system. The lack of a centralized, dependable system that makes it simple for people to find out

whether their prescription medications are available at different pharmacies and medical facilities increases these difficulties. Clinical services offered by pharmacists positively impact disease management, contributing to the broader spectrum of healthcare [2] [3]. The role of pharmacies, often the initial point of contact for patients, is pivotal in diminishing health disparities and strategically influencing patient health [4] [5] The lack of efficient ICT tools complicates the problem of pharmaceutical accessibility in a nation where the healthcare sector faces logistical and infrastructural limitations. Even if the Philippines has achieved progress in the field of healthcare, the lack of digital tools to find pharmacies offering the necessary medications makes it difficult to make the most use of the resources that are available.

Global Positioning System (GPS) and Geographic Information System (GIS) technologies play essential roles to the transformation of healthcare systems, particularly when it comes to medication location. By combining GIS and GPS, a potent toolkit can be developed for improving medication accessibility, resolving distribution issues, and improving public health outcomes [6] [7].

## II. RELATED LITERATURES

### A. Global Perspectives on Healthcare Accessibility

Access to Prescription medication in particular is a global issue with multiple factors to consider. In order to ensure every individual worldwide has access to necessary medical care without facing financial hardship; the World Health Organization (WHO) emphasizes the significance of Universal Health Coverage (UHC) [1]. Studies like the in [8], which examined the costs, accessibility, and availability of medications in 36 developing and middle-income nations and revealed structural and economic factors affecting pharmaceutical access, nevertheless, continue to show the persistence of worldwide inequities. Information technology is becoming a transformational force in the field of global health. The World Bank's "World Development Report 2021" notes how information technology, such as telemedicine and electronic health records, is enhancing access to and the quality of healthcare provided worldwide [9].

Tools for information and communication technology (ICT) have become vital for improving medicine delivery networks throughout the world. e-Prescribing platforms, digital health apps, and mobile health apps all help to provide better

patient access and effective prescription medication management [10]. Due to ongoing difficulties in pharmaceutical supply chains, organizations such as in [11] and [12] have launched the Access to Medicine Index, an assessment of pharmaceutical firms' global efforts to improve ethics and accessibility. The significance of inclusivity and innovation in tackling global healthcare accessible concerns is emphasized by collaborative worldwide efforts, patient-centric methods, and strategic collaborations [13].

### B. Factors Affecting Medication Accessibility in the Philippines

Medication accessibility in the Philippines has various problems influenced by the country's socioeconomic structure, healthcare system, and geographic dispersion. Affordability is a major concern, especially considering the wide economic gaps among the people. The Philippines has a mix of public and private healthcare services, and medicine out-of-pocket costs can be a burden for a lot of people, particularly those without appropriate insurance coverage [14]. High drug costs add to financial hurdles, influencing treatment adherence and may compromise health results.

The decentralization of the Philippines healthcare system and the geographic dispersal of healthcare institutions offer challenges to pharmaceutical accessibility, particularly in rural or isolated locations. In these areas, limited access to pharmacies and healthcare practitioners might result in increased travel costs and drug outages, adding to health inequities [15]. The government's initiatives to address these concerns, such as the passage of the Universal Healthcare Law, seek to enhance access to vital medications by reducing financial risk and extending healthcare coverage [16]. However, obstacles remain in ensuring that such policies are implemented effectively.

Disruptions in pharmaceutical delivery, for example, might have an influence on drug supply in the Philippines. The country has faced medicine shortages, impacting the population's access to some prescriptions [17]. Pharmaceutical accessibility in the Philippines is also influenced by regulatory constraints such as demanding licensing processes and limits on specific treatments. Understanding and addressing these characteristics is critical for devising tailored treatments to increase drug availability and adherence, especially given the country's unique healthcare system. Future research and policy initiatives should continue to investigate approaches to improve pharmaceutical accessibility while taking into account the distinct challenges and potential of the Philippine healthcare system.

Finally, healthcare delivery technologies such as telemedicine and digital health solutions have the potential to improve pharmaceutical accessibility by offering alternate channels for prescription distribution and monitoring [18]. Integrating these technologies into healthcare systems can increase convenience and accessibility, particularly for people with mobility issues or who live in rural places.

### C. Role of ICT in Pharmaceutical Services

The use of ICT in pharmaceutical supply chain management has substantially increased the distribution process's reliability and transparency. ICT also makes inventory management, order processing, and demand forecasting easier, ensuring that pharmaceutical items are accessible when and where they are required, resulting in a more dependable and responsive supply network. ICT-enabled pharmacy information systems and electronic health records (EHRs) have transformed patient care and drug administration. Telepharmacy services, a subset of ICT in pharmaceutical services, have been recognized, particularly in rural or disadvantaged regions, through offering patients with prescription consultations and professional guidance without the need to visit a physical pharmacy [19].

In the Philippines, Information and Communication Technology (ICT), specifically Geographic Information System (GIS) applications, are establishing themselves as a valuable tool for addressing difficulties linked to pharmaceutical supply and monitoring within the healthcare system. GIS technology proves essential for mapping and visualizing the geographical distribution of healthcare institutions, pharmacies, and pharmaceutical supply chains. This assists in identifying overlooked regions and improving the distribution of healthcare resources across the archipelago to promote more equal access to medications [20]. Furthermore, GIS applications serve to track and monitor drug availability by providing real-time data on the condition of pharmaceutical supplies in various locations, assisting in the management of shortages and improving overall supply chain management.

Integrating GIS into pharmacy information systems has the potential to transform pharmaceutical monitoring and prescription administration. GIS enables pharmacy geospatial mapping, allowing healthcare professionals to analyze the accessibility of pharmaceutical services and identify locations with drug distribution shortages [21]. This data is useful in establishing targeted actions to enhance medicine availability, particularly in rural or isolated areas. The use of GIS in pharmacy information systems can help improve medication adherence monitoring by providing insights into patient demographics and their proximity to healthcare institutions.

While incorporating ICT such as GIS into pharmaceutical services has significant advantages, problems such as data security, interoperability, and providing fair access to digital healthcare solutions must be carefully considered. To realize the potential of ICT in pharmaceutical services, ongoing research and deliberate implementation efforts are required, ensuring that technological innovations contribute favorably to the overall quality, safety, and accessibility of pharmaceutical care.

## III. METHODOLOGY

This study adopts the Rapid Application Development (RAD) approach, a methodological framework designed to address the drawbacks inherent in total system development methods [22]. The RAD model emphasizes on flexibility and adaptability facilitating the swift and cost-effective development of high-quality systems that can easily meet changing user needs. The method comprises four main phases: requirement planning, iterative development, system prototyping, and the throwaway prototype as shown in Fig. 1.

Central to this process is the initial development of the alpha version, where subsequent user testing and feedback inform the refinement of subsequent versions. This iterative process establishes a clear and linear understanding of the project's scope, enabling the development team to develop systems with extensive functionality within timelines [22].



Fig. 1.   Rapid application development.

### A. Planning

A web application and a mobile application are the two separate application models that are integrated into the system's architecture based on the requirements set by stakeholders. Every model has a different set of features and is designed to meet the needs of particular user groups based on discussions. While the mobile application is meant for general users and offers a user-friendly interface for a wider audience, the web application has been tailored for pharmacists to use, catering to their unique demands such as the management and inventory of medical drugs. The inherent features of mobile applications frequently contribute to their being considered as more adaptive which is more preferred by general users. Research on user-centric design concepts highlights how mobile app development must be customized to users' demands and touch interfaces and smaller screens [23]. Usability and responsiveness are often given first priority in these apps, which makes them ideal for mobile use. Furthermore, mobile applications can be easily integrated with device-specific functionality, such as cameras and GPS, to provide enhanced user experiences.

### B. User Design

The system comprises two integral components: the frontend and the backend as shown in Fig. 2. Together, these parts function collectively to provide an effective and logical experience.

The frontend, or user interface, is in the forefront and allows users to interact with the system. When users initiate requests, the frontend actively collects input data, creating a bridge between the user and the system. This input data is then encapsulated into JavaScript Object Notation (JSON), a lightweight data interchange format, preparing it to communicate with the backend.

The backend, which functions as the system's engine, receives the JSON data that has been processed. In this case, the complex task of processing the received request is handled by the backend, which makes use of a number of well-defined system logics to guarantee correct and effective handling. This covers tasks such as running algorithms, querying databases, and managing the system's general operation.

To ensure consistency in data transmission, the backend encodes the response into JSON format after processing is complete. The frontend receives this response, which presents the results of the user's request. The JSON respond to is then decoded by the frontend and formatted such that it can be readily read and understood by the user.



Fig. 2.   Pharmacy navigator architectural framework.

On the other hand, after three iteration and frequent consultations from stakeholders, two of whom were pharmacist, five general users, and two developers, features which are distinct and common to both applications are shown in Table I.

TABLE I.    LIST OF FEATURES, DISTINCT AND COMMON TO THE MOBILE AND WEB APPLICATION

| Application Model | Features |
|---|---|
| Web Application | The software is intended for use by pharmacists and their staff. It includes a thorough list of all the prescription medications that the pharmacy has on hand. A specialized component of the system is designed to manage the stock of pharmaceutical drugs that are available. This section offers a simplified interface to facilitate effective inventory management. |
| Mobile Application | Geographic information systems are incorporated into both models, which improves their functionality by offering location-based information like the closest pharmacy having the needed prescription medication.

Chat system is also incorporated to allow users to communicate with the pharmacy or the client vice versa. |
| | The application is exclusively designed for patients or general users. It incorporates a search feature that enables users to check the availability of specific pharmaceutical medicines. Upon initiating a search, the system provides information on the nearest pharmacy where the desired pharmaceutical drug is available. |

The system has been purposefully developed to meet the user requirements collected during the Rapid Planning Phase, translating these requirements into an effective system design. Table II and Table III provided below present a list of the functional requirements distinct for web and mobile applications respectively.

TABLE II.    LIST OF FUNCTIONALITY AVAILABLE FOR THE WEB APPLICATION

| Web Application | |
| --- | --- |
| *Functions* | *Actions* |
| **User Registration** | The system is expected to prompt an error message for unsuccessful registration and confirmation message if successfully registered. |
| **Login/Logout** | The system has a validation process for registered users.<br>If either or both of the username or password are/is invalid, the system is expected to prompt an error message.<br>The system will allow the entry of the user in the system of username and password are verified.<br><br>The system is expected to allow the user to logout from the system |
| **Display List of Medicines** | The system is expected to display all medicines stored in the pharmacy.<br><br>The system is expected to display the medicine searched by the user |
| **Manage Inventory** | The system will prompt a response message when updating, adding, and deleting inventory |
| **Manage Pharmacy Information** | The system allows the management of information and is expected to prompt a response message for every action performed. |
| **Chat System** | The system is capable of sending and receiving messages |

TABLE III.    LIST OF FUNCTIONALITY AVAILABLE FOR THE MOBILE APPLICATION

| Mobile Application | |
| --- | --- |
| *Functions* | *Actions* |
| **View List of Pharmacies** | The system is expected to show the list of pharmacies available in the database. |
| **Search Pharmacy** | The system is expected to display the information of the pharmacy, based on the information provided to the system. |
| **Get Pharmacy's Location** | This allows the user to get information of the pharmacy selected by tapping or navigating within the map. |
| **Get Pharmacy's Shortest Route** | The system is expected to show the shortest route leading to the selected pharmacy having the prescribed medicine. |
| **Search Medicine** | The system is expected to display the drug(s) information based on the search item. |
| **Display Medicine** | The system is expected to display all stored medicine in the database. |
| **Chat System** | The system is capable of sending and receiving messages |

## C. Testing

For the usability test of the mobile application, the study targeted random participants aged eighteen (18) and above who met at least one of the following criteria: (a) residing in Iligan City, (b) having studied in Iligan City, or (c) currently enrolled in a school in Iligan City.

This research study includes the voluntary participation of respondents with basic understanding of pharmaceuticals and be able to perform activities such as naming, classification, and categorization of pharmaceutical goods. Participants are not required to live or study in Iligan City.

Demographic information was collected as part of the questionnaire. To uphold ethical standards, letters of consent were provided to participating respondents. A total of forty-five respondents have tested the mobile application while there were only five pharmacists for the web application.

The tests were conducted online, utilizing the "Play Store" for the mobile application (for general users) and a URL for the web application (pharmacist/staff role). The main goals of the alpha test, which is the initial assessment stage, were to get relevant user feedback and validate that the developers had successfully complied with the user system requirements.

To begin the testing procedure, participants were given a Google Form with instructions and tasks. This method allowed for "freestyle" testing, allowing participants to explore the application in their own way. The given activities were designed to determine if the application efficiently meets the research questions and preserves the integrity of its functionality, especially in the context of the mobile application and web application. A survey was given after the task. This survey was essential in determining how well the application met the users' expectations and requirements. It examined a number of issues, such as overall application satisfaction, interface design, and user experience. A five-level Likert scale was employed for rating each question.

A key part of the survey was obtaining qualitative input from users. It was encouraged for participants to share thoughts and recommendations, especially with reference to areas where changes or enhancements could improve their experience. Users were encouraged to provide their opinions on possible improvements by answering questions like "Is there anything else you would like us to improve or change?" in order to gather insights that might not be addressed by initial activities.

## IV. RESULTS AND DISCUSSION

The design, functionality, and user interfaces of the desktop and mobile applications are assessed and compared in this section. Aiming to improve stakeholders' engagement and satisfaction, developers designed at incorporating all of their requirements in the digital interfaces and functionalities.

Incorporating participants' feedback from the alpha testing is enhanced in the beta testing phase, which is the final stage of assessing the mobile application in the research study. Coherence with the parameters and constraints of the research is certain through this method. During the alpha testing, the researchers identified human errors, which prompted a series of UI/UX re-engineering for the program.

## A. Mobile Application

This mobile application is intended for patients or general users, and it was developed to a user-friendly search option that lets people locate whether certain prescription medications are available. The following are the features of the mobile application.

*1) View list of pharmacies, search pharmacy and get pharmacy's location:* This functionality enables users to navigate, tap, and untap specific pharmacy using the participant's geo-coordinates, facilitating interaction with the

application map. A persistent bottom sheet emerges on the bottom left, displaying the pharmacy's name and a "view" button that directs the user to the pharmacy's dedicated page.



Fig. 3. Screenshot of the get pharmacy's location feature in the alpha test.

The challenge with this interface as shown in Fig. 3 during the alpha test was the number of underlying pages, which could be minimized for a clearly defined functionality. As a result, it becomes diverted by the multitude of page elements and deviates from its primary goal of locating a pharmacy.



Fig. 4. Screenshot of the get pharmacy's location feature in the beta test.

Fig. 4 displays an improved iteration of the Get Pharmacy's Location Feature, incorporating feedback obtained during the alpha test. In this version, the search function is exclusively in the persistent bottom sheet, minimizes pages and buttons. Upon conducting a search for a specific medicine, the system displays all pharmacies offering the desired drug, allowing users to tap on a location/pharmacy to access detailed information. Each pharmacy's location is displayed, allowing users to make informed decisions based on proximity and

convenience. The objective is to provide a user-friendly interface that promotes ease of use and accessibility.

*2) Get pharmacy's shortest route:* An essential component of the application's testing phase, this functionality uses the Google Map API in order to guide visitors from their starting place to the selected pharmacy. Based on the user's starting location and preferred method of transportation (walking, cycling, or driving), participants assess how well the program performs in providing directions. Walking is chosen as the default method to consider the means of transportation by the majority of users. Furthermore, the feature provides users with several route possibilities, enabling them to select the most appropriate direction. In order to enhance the entire process of decision-making, this user-centric approach also includes crucial information like arrival time, distance, and route means.

Similar with the Get Pharmacy location, the alpha test result pointed out that this feature loses efficiency because it takes up numerous pages rather than being contained to one, making it difficult to use as shown in Fig. 5.



Fig. 5. Screenshot of the get pharmacy's shortest route in the alpha test.

With the influence from Google Maps for navigation, showed that a single page approach was more practical and user-friendly which is shown in Fig. 6. In addition to being in line with this preference, combining the two pages into one, improves interaction and simplifies the functionality for users.

*3) Search medicine and display medicine:* This tool makes it simple for consumers to search for a specific medication. It executes a case-insensitive search algorithm by interpreting the user's input and taking into account both the brand and generic names of medications. A list of medications that match the keyword entered appears as the result, making the search process simpler and user-friendly.

Fig. 6.   Screenshot of the get pharmacy's shortest route in the beta test.

The researchers discovered that one issue with the previous version of the feature was that it displayed extraneous information and was not engaging enough as shown in Fig. 7. Instead of needing to present the results on another page, smart predictions will appear below the search box and, when clicked, will divert to another website.



Fig. 7.   Screenshot of the search medicine in the alpha test.

In order to enhance user convenience, when a user chooses their preferred medicine, pertinent medical information will instantly display within the same window. This design improvement attempts to provide consumers fast access to important facts, reducing the need for further clicks and producing a more streamlined and user-friendly interface as shown in Fig. 8.



Fig. 8.   Screenshot of the search medicine in the beta test.

### B.  Web Application

The web application's core functionalities; Registration, Login, Medicine Display, Inventory Update, and Pharmacy Management, performed as expected during usability testing, aligning with the positive ratings provided by respondents. While the overall functionality met expectations, constructive feedback from one participant emphasized the need for enhancement in the UI/UX design. This feedback emphasizes the need of refining the user interface's perceived plainness and lack of vitality.

Despite the feedback, the Pharmacy Management tool was a significant feature. This module is excellent for developing trust in both the general public and pharmacy management. The pharmacy must register the pertinent data in the Register Pharmacy feature, as depicted in Fig. 9 to keep track of important pharmacy information. It works by including a pharmacy log, which carefully documents every activity taken by a certain pharmacy. This thorough documentation process maintains transparency and prevents pharmaceutical brands from being swapped. As a result, this feature adds greatly to the pharmacy's operational integrity.

The Display Medicine functionality, shown in Fig. 10, which allows the system to display every medication stocked in the pharmacy, is strength of the program. Notably, one respondent suggested that the presentation be improved by including the medication expiration dates for tracking. In addition, another respondent recommended changing the "Availability" label to "Stock" for consistency. This recommendation would include existing pharmaceutical information such as Generic Name, Brand Name, Dosage, Form, and Price, as well as real-time stock levels. The refinement intends to provide management a picture of the pharmacy's inventory, supporting informed decision-making and effective pharmaceutical resource management.

Fig. 9.   Screenshot of the register pharmacy in the beta test.



Fig. 10. Screenshot of the display / manage medicine in the alpha test.

One of the distinguished features of the system is the incorporation of a messaging system that allows users to communicate directly with any pharmacy's representative, enabling responsive and customized engagement environment. The feature was designed to serve as a dynamic link for users for inquiries, seek information, or address any concerns. As the system develop from the alpha test phase to the more polished beta version, the focus of development shifts to the complexities of its user interface. This shift is more than just a standard upgrade; it is a concerted attempt to improve the user experience. The refining process is extremely responsive to user feedback, ensuring that the system not only meets the different demands of its user base, but also delivers an intuitive platform.

### C.  For Both System

Fig. 11 shows the interface of the chat system incorporated into the system. Users can navigate questions about dose, possible side effects, and other options through interactions, which promotes a more informed and decision-making process. To put it simply, the integration of this chat system simplifies communication while also providing people with a useful tool for getting precise and customized information on their medical needs and health.

### D.  Testing and Evaluation

A majority found the implementation of the Google Map API, with the utilization of geo coordinates, to be successful in its early stages. A substantial 37 respondents, as shown in Fig. 12 were extremely satisfied with the functionality of the Search and Locate Pharmacy. It is clear that the majority of respondents thought the feature was useful and functional, despite the fact that one and five respondents, respectively,

expressed a modest satisfaction and neutrality with the feature. The majority of participants' overall favorable response highlights the usefulness of the Search and Locate Pharmacy capability, demonstrating that it satisfied their needs and expectations.

Forty-four respondents were satisfied with the functionality of the "View Available Medicine/Drug" feature as shown in Fig. 13. One of the comments emphasized how simple it was to locate particular medications, indicating that the search medicine feature answers the needs of most users. This feedback means how well the functionality works to give consumers a practical and easy-to-use experience when looking through the medications that are accessible.

A word cloud of comments received during the alpha testing stage is shown in Fig. 14. Respondents generally expressed satisfaction for the application's overall convenience, ease of use, and user-friendliness. Users responded positively to the application's design and use, highlighting its simple navigation and easy operations.  A subset of feedback highlighted concerns related to the redundancy of pages and experience (UI/UX) during testing. In response, the research team proactively addressed these issues, focusing on the refinement of redundant pages through a comprehensive re-engineering process.



Fig. 11. Screenshot of the chat system.



Fig. 12. Result of the locate and navigate pharmacy.

Fig. 13. Result of the view available medicine / drug.



Fig. 14. Word cloud of the application during the alpha test.

The word cloud in Fig. 15 depicts the input gathered throughout the beta testing stage. Remarkably, more users acknowledged the application as an improved version over the alpha test. The redesign prioritized resolving issues brought up in previous testing, making the experience more efficient and user-friendly. In particular, unnecessary pages and intricate content were improved for more basic versions, in line with user expectations for a more refined application.



Fig. 15. Word cloud of the application during the beta test.

In addition to verifying the changes made, the beta testing stage reaffirmed at improving the application in response to user feedback. The purpose is to make sure that the application satisfies the needs and preferences of a wide range of users by optimizing the user experience and improving the way information is presented.

## V. CONCLUSION

In conclusion, the development and implementation of the web and mobile application integrating Geographical Information System (GIS) for locating prescribed medicine have been driven to enhance medication accessibility, distribution challenges, and improve public health outcomes [24] [25]. The adoption of the Rapid Application Development (RAD) methodology, having continuous engagement with stakeholders aims to meet user needs and requirements in efficiently locating prescribed medicines from nearby pharmacies. The mobile application, equipped with features such as viewing medicine lists, searching pharmacies, obtaining pharmacy locations and shortest routes, searching for specific medicines, displaying medicine details, and incorporating a chat system, serves as a user-friendly platform designed for the general public in managing their healthcare needs in terms of searching prescribed medicine. Simultaneously, the web system offers features like displaying medicine lists, managing inventory, overseeing pharmacy information, and integrating a chat system to cater to a wider range of user needs. The positive feedback received during both the alpha and beta testing phases further validates the success of the application's functionalities. Respondents consistently reported that the features were easy to use, easy to navigate, and aligned with their needs. This user-centric approach, with the set of features, positions the integrated web and mobile application as a valuable resource in locating medicines, and contributing to the overall improvement of public health outcomes.

## REFERENCES

[1] World Health Organization (WHO), "Tracking Universal Health Coverage: 2019 Global Monitoring Report.," [Online]. Available: https://www.who.int/healthinfo/universal_health_coverage/report/2019/en/. [Accessed January 2023].

[2] M. d. C. V. B. Sousa, B. D. Fernandes, A. A. Foppa, P. H. R. F. Almeida, S. d. A. M. Mendonça and C. Chemello, "Tools to prioritize outpatients for pharmaceutical service: A scoping review," *Research in Social and Administrative Pharmacy,* vol. 16, no. 12, pp. 1645-1657, 2020.

[3] B. D. Fernandes, P. H. R. F. Almeida, A. A. Foppa, C. T. Sousa, L. R. Ayres and C. Chemello, "harmacist-led medication reconciliation at patient discharge: A scoping review," *Research in Social and Administrative Pharmacy,* vol. 16, no. 5, pp. 605-613, 2020.

[4] K. Hurley-Kim, J. Unonu, C. Wisseh, C. Cadiz, E. Knox, A. Ozaki and A. Chan, "Health Disparities in Pharmacy Practice Within the Community: Let's Brainstorm for Solutions," *Front Public Health,* vol. 8;10:847696, 2022.

[5] M. L. Ilardo and A. Speciale, "The Community Pharmacist: Perceived Barriers and Patient-Centered Care Communication," *Int J Environ Res Public Health,* vol. 17(2): 536, no. 10.3390/ijerph17020536, 2020.

[6] B. K. H. T. J. V. L. J. A. S.-P. A. R. Neda Firouraghi, M. Furs, L. Salvador-Carulla and N. Bagheri, "The role of geographic information system and global positioning system in dementia care and research: a scoping review," *Int J Health Geogr,* vol. 21, no. 8, pp. https://doi.org/10.1186/s12942-022-00308-1, 2022.

[7] G. Musa, P.-H. Chiang, T. Sylk, R. Bavley, W. Keating, B. Lakew, H.-C. Tsou and C. Hoven, "Use of GIS Mapping as a Public Health Tool–-From Cholera to Cancer," *Health and Nursing,* p. https://doi.org/10.4137/HSI.S10, 2013.

[8] A. Cameron, M. Ewen, D. Ross-Degnan and D. Ball, "Medicine Prices, Availability, and Affordability in 36 Developing and Middle-Income Countries: A Secondary Analysis," *Lancet,* Vols. 373(9659):240-9, 2009.

[9] The World Bank, "DATA FOR BETTER LIVES," International Bank for Reconstruction and Development, Washington DC, 2021.

[10] P. De and M. Pradhan, "Effectiveness of mobile technology and utilization of maternal and neonatal healthcare in low and middle-income countries (LMICs): a systematic review.," *BMC Women's Health,* Vols. 23, 664 , pp. https://doi.org/10.1186/s12905-023-02825-y, 2023.

[11] Access to Medicine Foundation, "Access to Medicine Index," The Netherlands, 2021.

[12] P. Yadav, "Health Product Supply Chains in Developing Countries: Diagnosis of the Root Causes of Underperformance and an Agenda for Reform," *Health Systems & Reform,* vol. 1:2, pp. 142-154, 2015.

[13] S. Flessa and C. Huebner, "Innovations in Health Care—A Conceptual Framework," *Int J Environ Res Public Health,* vol. 18(19): 10026., 2021.

[14] X. Javier, P. Crosby, R. Ross, M. E. Ranchez-Vila and M. S. Santos, "Understanding Out-of-Pocket Expenditure for Outpatient and Inpatient Care," 2022. [Online]. Available: http://www.healthpolicyplus.com/ ns/pubs/18653-19126_PhilippinesOOP.pdf. [Accessed October 2023].

[15] M. J. N. Naria-Maritana, G. R. Borlongan, M.-A. M. Zarsuelo, A. K. G. Buan, F. K. A. Nuestro, J. A. Dela Rosa, M. E. C. Silva, M. A. F. Mendoza and L. R. Estacio, "Addressing Primary Care Inequities in Underserved Areas of the Philippines: A Review," *Acta Medica Philippina,* vol. 54, no. 6, pp. 722-733, 2020.

[16] A. M. Amit, V. C. Pepito and M. Dayrit, "Advancing Universal Health Coverage in the Philippines through self-care interventions," *The Lancet Regional Health - Western Pacific,* p. https://doi.org/10.1016/j.lanwpc. 2022.100579, 2022.

[17] "Presidential Decree No. 651, s. 1975," Official Gazette, [Online]. Available: https://www.officialgazette.gov.ph/1975/01/31/presidential-decree-no-651-s-1975/. [Accessed 10 June 2023].

[18] A. Haleem, M. Javaid, R. P. Singh and R. Sumanc, "Telemedicine for healthcare: Capabilities, features, barriers, and applications," *National Center for Biotechnology Information,* vol. 2:100117, 2021.

[19] S. Baldoni, F. Amenta and G. Ricci, "Telepharmacy Services: Present Status and Future Perspectives: A Review," *Medicina (Kaunas),* vol. 55(7):327, 2019.

[20] C. T. Jagadeesan and V. J. Wirtz, "GIS technology proves essential for mapping and visualizing the geographical distribution of healthcare institutions, pharmacies, and pharmaceutical supply chains. This assists in identifying overlooked regions and improving the distribution of healthcare," *Journal of Pharmaceutical Policy and Practice,* 2021.

[21] B. D. Fernandes, A. A. Foppa, P. H. R. F. Almeida, A. Lakhani and T. d. M. Lima, "Application and utility of geographic information systems in pharmacy specific health research: A scoping review," *Research in Social and Administrative Pharmacy,* vol. 18, no. 8, pp. 3263-3271, 2022.

[22] R. Delima, H. B. Santosa and J. Purwadi, "Development of Dutatani Website Using Rapid Application Development," International Journal of Information Technology and Electrical Engineering, vol. 1, no. 2, p. https://doi.org/10.22146/ijitee.28362 , 2017.

[23] L. Punchoojit and N. Hongwarittorrn, "Usability Studies on Mobile User Interface Design Patterns: A Systematic Literature Review," *Advances in Human-Computer Interaction,* vol. 2017, p. https://doi.org/10.1155 /2017/6787504.

[24] C. Smith and J. Mennis, "Incorporating Geographic Information Science and Technology in Response to the COVID-19 Pandemic," *Preventing Chronic Disease,* vol. 9;17:E58, 2020.

[25] C. Vîlcea and S. Avram, "Using GIS methods to analyse the spatial distribution and public accessibility of pharmacies in Craiova city, Romania," *Bulletin of Geography Socio-economic series,* vol. 45(45), pp. 125-132, 2019.

# Hybrid Bio-Inspired Optimization-based Cloud Resource Demand Prediction using Improved Support Vector Machine

Nisha Sanjay[1], Dr. Sasikumaran Sreedharan[2]

Research Scholar[1], Research Supervisor[2],

Faculty of Computer Science and Multimedia, Lincoln University College Marian Research Centre,

Marian College Kuttikanam (Autonomous), Kerala, India[1, 2]

*Abstract*—In order to furnish diverse resource requirements in cloud computing, numerous resources are integrated into a data centre. How to deliver resources in a timely and accurate manner to meet user expectations is a significant concern. However, the resource demands of users fluctuate greatly and frequently change regularly. It's possible that the resource provision won't happen on time. Furthermore, because some physical resources are shut down to save energy, there may occasionally not be enough of them to meet user requests. Therefore, it's critical to offer resource provision proactively to ensure positive user involvement using cloud computing. To enable resource provision in advance, it is essential to accurately estimate future resource demands. Using machine learning techniques, we offer a unique approach in this study that tries to identify key features, accelerating the forecast of cloud resource consumption. Finding the classification method with the greatest fit and maximum classification accuracy is crucial when predicting cloud resource consumption. The attribute selection method is used to decrease the dataset. The categorization process is then given the reduced data. The hybrid attribute selection method used in the investigation, which combines the bio-inspired algorithm genetic algorithm, the pulse-coupled neural network, and the particle swarm optimization algorithm, improves classification accuracy. The accuracy of prediction employing this technique is examined using a variety of performance criteria. When it comes to predicting the demand for cloud resources, the experimental results show that the suggested machine learning method performs more effectively than traditional machine learning models.

*Keywords*—*Cloud computing; resource demand; machine learning; cloud resource demand prediction; bio-inspired algorithm*

## I. INTRODUCTION

With the benefits of characteristic features like resource pool and a pay-as-you-go paradigm, cloud computing has found widespread use in a variety of industries. Infrastructure as a Service (IaaS) is a brand-new cloud service paradigm which offers clients virtual machines as resources (VMs). Accurate and timely allocation of VMs to client tasks in IaaS service model is a major concern. Some cloud providers continue to offer VMs statically, which results in increased operational cost for customers and reduced resource usage for cloud providers. A better alternative is to create VMs dynamically in response to the current resource demands.

However, resource demands fluctuate significantly at times and change continuously throughout time. Users can quickly apply for numerous VMs, for instance. This causes these VMs to take a very lengthy time to create. Due to some of the physical resources being shut down to save energy, even those that are currently in operation may not be enough to meet user requests. For this reason, developing a proactive resource provision is essential to guaranteeing that customers get a positive cloud computing experience. In response to anticipated resource demands, proactive resource provision might offer resources in advance. However, if overestimated, this can be a resource waster.

It goes without saying that if the resource demand is estimated to be less than the actual requirements, user demand of the resources are not met. Therefore, the main challenge is to effectively forecast future resource demands to reduce overestimation and underestimation. An overview of the challenges and approaches for forecasting usage of resources in cloud computing can be found in the study [1]. To prevent resource over-provisioning, Chen et al. [2] developed a forecast method exclusively for burst workload. To eliminate bursts and sounds, this technique employs the Fast Fourier Transform (FFT) algorithm, which increases prediction accuracy. A cloud workload prediction model is put out by Roy et al. [3]. This model predicts future workload using autoregressive moving average method of the second order (ARMA) and then uses a performance model of the average app response time to forecast resource requirements.

Using ensemble models, two self-adaptive resource demand prediction techniques are proposed [4, 5]. In order to increase prediction accuracy, Xu et al. [6] propose the GFSS-ANFIS/SARIMA prediction model, which integrates Seasonal Autoregressive Integrated Moving Average Model (ARIMA) and Generalized Fuzzy Soft Sets with Adaptive Neuro Fuzzy Inference System. Data mining and statistical techniques are used in Verma et al. [7] resource prediction framework for multi-tenant service clouds to forecast resource demands in order to minimize resources and provisioning time. To obtain precise performance forecasts in hybrid clouds, Imai et al. [8] suggest a model which has workload-tailored elastic compute unit (WECU) as a computing power unit. Brown's quadratic exponential smoothing approach is used by Mi et al. [9] to

forecast activities ahead and adjust resources as needed in a data center.

The resource prediction method used by Minarolli et al. [10] incorporates the cross relation of asset utilization among VMs that are part of the same application. Bankole et al. [11] proposed three performance prediction models that employ neural networks (NN), linear regression (LR), and support vector regression (SVR) methods independently. The experimental findings demonstrate that SVR is the favored model, outperforming other models in terms of prediction accuracy. Similar to this, certain machine learning techniques, such as the multi-layer perceptron (MLP) [12], Support vector regression (SVR) [12], deep belief network (DBN) [13], and artificial neural network (ANN) [14], are used to forecast resource utilizations. There are two categories of prediction. The former involves statistical methods and the latter employs machine learning techniques. Even though machine learning methods can produce predictions with a higher degree of accuracy, it requires setting up the training model and extracting the features from a humungous amount of data. Though sampling data might occasionally change greatly and be insufficient as inputs for machine learning techniques. The prediction accuracy of statistical methods is often poor for non-stationary and non-linear data.

One such well-known metaheuristic method based on swarm-intelligence is particle swarm optimization (PSO), which has demonstrated its superiority in resolving a variety of real-world optimization issues from fluid mechanics, wireless sensor networks, engineering, applied sciences and academia [15]. Despite its success in locating competitive solutions, the PSO still has trouble sustaining strong exploration and is susceptible to becoming stuck in local optima, which leads to an insufficient exploration-exploitation tradeoff [16]. These flaws have an impact on the subsequent quality of the scheduling solution. Additionally, the "No Free Lunch Theorem (NFL)" [17] establishes the impossibility of finding a single metaheuristic algorithm capable of efficiently solving all optimization issues. These facts serve as powerful impetuses for the current research project, which proposes a hybrid PSO-based scheduling solution to get over the constraints of regular PSO by combining it with genetic algorithm (GA) and pulse coupled neural networks (PCNN).

The remainder of the paper is organized as below. Section II includes an overview of prior studies in the same field. The methodology, system workflow, feature extraction, and demand prediction employed in the proposed approach are covered in Section III. Section IV includes the performance analysis of the proposed hybrid solution. Section V provides a conclusion, marking the end of the paper.

## II. RELATED WORK

To solve cloud based scheduling (CBS) difficulties, current research uses heuristics and metaheuristics, such as shortest job first (SJF), First come first served (FCFS), Max-min, Min-min, Minimum completion time (MCT), Minimum execution time (MET), and Suffrage [18]. Heuristics offer problem-specific solutions and are apt for solving minor problems. In contrast, ensemble methods (MHs), are simple, iterative, adaptable, highly speculative, algorithms that direct a subordinated

heuristic through smart mechanism [19]. For complicated and larger scheduling issues, metaheuristics-based scheduling solutions have outperformed problem-specific heuristics [20]. However, MHs typically have a few flaws, such as premature convergence, getting stuck in neighboring best value, an absence diversity, and imbalance among the examination and development stages of the energy spectrum [21]. If these flaws are applied to work scheduling issues, the results can be undesirable.

The research has also advocated the use of combination heuristics to address the drawbacks of solo metaheuristics [22]. Several whale optimization algorithms (WOA)-based scheduling approaches have already been put out for scheduling bag-of tasks (BoT) applications to get results that are almost optimal and are motivated by the humpback whales' hunting method. These solutions include those that employ conventional, customized, and fusion of WOA techniques [23]. A standard WOA, Gaussian model, and opposition-based learning (OBL) approaches are combined and used in a cloud scheduling solution called GCWOAS2 [24] to provide effective task-resource couples. One more current study paper [23] proposes a combination metaheuristic approach dubbed OWPSO to address the shortcomings of the original WOA by combining OBL and PSO algorithms. The authors of [25] proposed random double adaptive WOA (RDWOA) employing the Bee optimization algorithm techniques for arranging cloud workloads to decrease implementation cost and time.

The increase and decrease operators might be augmented to the typical WOA to enhance search efficiency, according to authors in [26], who also proposed utilizing two advanced optimization algorithms. In study [27], authors proposed an Improved WOA for Cloud task scheduling (IWC) algorithm that makes use of the local weight method to enhance neighboring explore effectiveness and prevent the basic WOA's early convergence. WOA and harmony search algorithm (HS) were hybridized to create WHOA by the authors of [28] in order to reduce execution costs and energy usage. The WOA-based cloud task scheduling solution, which simultaneously optimizes makepan and energy usage, was proposed by Sharma and Garg [29]. To schedule BoT applications over clouds, the whale-Scheduler technique was recommended, yielding the best makespan and execution cost [30]. There have been many different GA-based scheduling approaches proposed in the past, including the basic GA method [31], improved GA strategies employing modified mutation as well as crossover procedures, and fusion of GA results that combine classic GA among additional methodologies [32].

Numerous studies have greatly improved scheduling efficiency over cutting-edge heuristics by using basic variants of SOS and particular improvements to Symbiotic Organism Search (SOS) algorithms utilizing chaotic sequences with dissent learning [33]. To overcome CBS concerns and achieve different QoS objectives, researchers have suggested conventional ant colony optimization (ACO) as well as other altered ACO-based optimization techniques [34]. In numerous researches, classic PSO algorithms, modified PSO variations, and hybrid PSO variants have all been applied to handle CBS

situations in order to achieve a variety of objectives, including minimizing computation time and complexity and addressing load balance challenges [35]. An adjustable inertia weighting method is suggested for the CBS problem to tradeoff between exploration and exploitation [36]. For maintaining population variety and enhancing solution quality, Chen and Long [37] proposed a fusion of optimization approach integrating PSO and ACO algorithms. A bi-objective PSO scheduler was used by the developers of [38] to improve system performance and lower execution costs. A multi-objective PSO strategy is used in two deadline-constrained scheduling techniques to enhance QoS metrics values [39].

A non - linear and non PSO was employed by the authors of Ref. [40] to lessen the time for scheduling the workload. Several strategies for workload allocation in a cloud have been suggested, using both a standard Cuckoo Search (CS) implementation [41] as well as modifications to the CS's basic structure and the incorporation of additional metaheuristics. Chhabra et al. proposed a fusion of CS method in [42] to enhance the exploration potential of standard CS and to achieve more appropriate scheduling than current generation heuristics. This metheuristic merged CS and DE algorithms. a solid metaheuristic in the form. GWO has been employed in the past to create almost ideal scheduling solutions to improve various QoS metrics. To optimize both makespan and energy, for instance, a MO-GWO technique was recommended [43]. Modified GWO is formulated in [44] with alterations on the fitness function to consider makespan and cost. In [45], mean GWO is studied to achieve better performance tackling scheduling concerns. It results in lower makespan and reduced energy usage.

Elaziz et al. in [46] integrated the DE algorithm's effective local searching feature and Moth search (MS) method for better scheduling solution. The Bacterial Foraging Optimization (BFO)-based scheduling strategy has been recommended by Milan et al. [47] to optimize extent of imbalance, idle time, and overall execution time. The study in [48] proposes a novel approach to cloud computing scheduling problem solution: Water Pressure Change Optimization (WPCO). The phenomena of water density changing as pressure is increased due to changes in the physical properties of water serves as the inspiration for the novel WPCO technology. WPCO provides the best solution quality in comparison to the standard metaheuristics. According to the authors of [49], a scheduling strategy based on social group optimization algorithm (SGO) is proposed for a diverse cloud environment that can be used to resolve CBS issues with the highest possible throughput and the shortest possible makespan.

Inadequate exploration and utilization process balancing, slow convergence, a failure to focus on schedule order optimization, a lack of products developed using standard workloads, a lack of numerical solutions to tune metaheuristic variables, and concurrent performance and energy consumption optimization are common problems or limitations of the current research studies based on metaheuristic approach for scheduling BoT applications over cloud systems. These flaws leave a lot of room for developing new metaheuristics or

refining already-existing ones to increase the CBS problem's efficiency.

## III. PROPOSED APPROACH

Fig. 1 depicts the planned work's entire organizational structure. It consists of two phases namely training and using the constructed model for prediction. The major goal of the training phase is to learn about the cloud resource requirements and how to predict resource demand from provided data sets. As a result, it is referred to as the training phase. At this point, a method called attribute selection is used to cut down on the amount of data that was collected. The correlating attributes are identified when the data set's dimensions are reduced. These characteristics are crucial for forecasting cloud resource consumption.



Fig. 1. Organizational structure of the proposed approach.

The steps of training and prediction are the same. The similar procedure was performed throughout prediction for the test data. Finally, the anticipated outcomes will be attained. In the following section, each procedure is thoroughly explained. The following list of modules contains the working stages of the planned work.

1) Attribute Selection using PCGPSONN.
2) Correlation feature extraction.
3) Cloud Resource Demand Prediction using SVM.

### A. Attribute Selection using PCGPSONN

Only the unique Pulse Coupled Genetic Particle Swarm Optimization Neural Network (PCGPSONN) is used to determine the cloud resource database's most crucial properties. The attributes must be carefully chosen for an accurate demand prediction of cloud resources. Low accuracy, prediction inaccuracy, or failure might result from the incorrect selection of these attributes. If the feature selection initial selection was incorrect, the approach will never reach the global minimum and more runs will only take the algorithm to a local minimum. The best position for each particle, Pg, as determined by the Genetic Particle Swarm Optimization (GPSO) algorithm, is searched for in this work using the Pulse Coupled Neural Network (PCNN) algorithm.

In GPSO method, the GA operators and the PSO update mechanism typically operate with the same population throughout initialization step. Uniformly distributed random numbers should be used to create the initial population of all its members. Therefore, the attribute selection process must go through more iterations because of this random distribution's sluggish convergence. However, in our suggested method, GPCNN solutions are used to allocate the PSO's initial population. GPCNN and PSO split the entire number of iterations evenly. GPCNN runs the first half of the iterations, and the answers are provided as the PSO's initial population. PSO is in charge of the final iterations. Thus, the issue of delayed convergence is resolved, and the attribute selection process requires fewer iterations.

Additionally, the hybrid approach we suggest should include local and worldwide search. Then, for each particle, generate the best position Pg, and the attribute selection method is then given these best positions Pg to further refine the search process. Consequently, our hybrid strategy performs better than this approach. The PCGPSONN algorithm is displayed below.

**Algorithm of PCGPSONN**:

Input:

- Attributes (X) and their values.
- Population Size (pop_size): Number of individuals in the population.
- Max Generations (max_gen): Maximum number of generations.
- Crossover Probability (crossover_prob): Probability of crossover occurring.
- Mutation Probability (mutation_prob): Probability of mutation occurring.
- Particle Swarm Size (swarm_size): Number of particles in the swarm.
- Inertia Weight (inertia_wt): Weight for the inertia term in PSO.
- C1 and C2 (c1, c2): Constants for the cognitive and social components in PSO.

Output:

- Best Attribute (S): The best solution found.
- Fitness Value (W): The fitness value corresponding to the best solution.

1. Initialize Population:

   - Generate an initial population of solutions (pop_size) randomly.

2. Loop through Generations:

   - Repeat for a specified number of generations (max_gen) or until a convergence criterion is met.

   2.1 Evaluate Fitness:

   - Evaluate the fitness W of each solution X in the population.

2.2 Genetic Algorithm Steps:

- Select solutions S1 for crossover based on their fitness.
- Perform crossover with a certain probability (crossover_prob).
- Mutate selected solutions with a certain probability (mutation_prob).

2.3 Particle Swarm Optimization (PSO) Steps:

- Initialize a particle swarm (swarm_size) with positions and velocities.
- Evaluate the fitness W of each particle X.
- Update particle positions and velocities based on PSO equations.
- Track the global best position S2 found by the swarm.

2.4 Combine Genetic and PSO Steps:

- Replace the worst solutions S1 in the population with the best particles from the swarm S2.

2.5 Update Population:

- Create a new population for the next generation by combining the modified population and the PSO swarm.

3. Return Best Solution:

   - Return the best solution S found in the final population.

4. Execute PCNN:

- Execute PCNN with W and S

PCNN algorithm

Alpha_F = 0.1 Decay term for feeding

Alpha_L = 1.0 Decay term for linking

Alpha_T = 1.0 Decay term for threshold

V_F = 0.5 Magnitude scaling term for feeding

V_L = 0.2 Magnitude scaling term for linking

V_T = 20.0 Magnitude scaling term for linking

Beta = 0.1 Linking strength

Num = 100 The number of iterations

W= [0.5 1 0.5;1 0 1;0.5 1 0.5] Initial values for W

M = [0.5 1 0.5;1 0 1;0.5 1 0.5] Initial values for M

F = zeros(size(S)) Initial values for F

L = F Initial values for L

Y = F Initial values for Y

U = F Initial values for U

T = Ones(size(S)) Initial values for T

S = im2double(S) Normalizing to lie within [0,1]

for n = 1:Num

F = exp(_Alpha_F) *F + V_F*conv2(Y,W) + S Update the feeding input

L = exp(_Alpha_L) * L + V_L*conv2(Y,M) Update the linking input

U = F_ *(1 + Beta*L) Compute the internal activation

Y = double(UiT) Update the output

T = exp(_Alpha_T) *T + V_T*Y Update the threshold input

End

The genetic algorithm's fitness value is W in this case and S is the best quality of genetic algorithms. The proposed approach chooses seven of the 11 attributes that are present in the cloud resource dataset. Table I provides a description of the complete attributes. Table II provides descriptions of the chosen features.

Based on the values of the attribute weights, PCGPSONN chooses the attribute. Instead of using binary presentation, a population of 200 clients used real-valued representation since the parameter coefficients were expressed using real-valued numbers rather than just 0 and 1. Seven separate attribute weight sets totaling 28 qualities made up each individual. The beginning population's members were chosen using machine learning techniques and expert-set weights. More precisely specified is the initial population. The PCGPSONN used a uniform crossover with discrete recombination for offspring creation and a roulette-wheel selection for parent selection. 80.0% of the time was spent on the crossover, and each gene's crossover points were selected separately and arbitrarily. The gene underwent mutation with a likelihood of 1.0% and was uniformly carried out by selecting a random value at random from the range and setting it as the new value at the present place. Additionally, elitism was employed during runs to preserve the greatest person among the population. The weight set that performed the best during evolution was something we did not want to lose. If the total population was higher than 21 at the culmination of the generation, a strategy for survivor screening was used. Those with the lowest categorization accuracy were removed from the population after the individuals were graded according to their accuracy. After 20 generations, the PCGPSONN algorithm concluded, or sooner if the best classification accuracy remained constant during a period of 10 iterations. Additionally, the examination came to an end if every member of the population were the same. The proposed PCGPSONN approach is contrasted with the GPSO and GPCNN approaches to demonstrate its efficacy. Tables III and IV display the GPSO and GPCNN results, respectively.

The selected attribute list of GPSO is shown in Table III. In GPSO approach first genetic algorithm is completed to get the fitness values of all attributes and then it is given to the PSO approach to complete the attribute selection process.

TABLE I.  FEATURES OF CLOUD RESOURCE DATASET

| S.No | Feature name |
|---|---|
| 1 | Timestamp |
| 2 | Disk read throughput |
| 3 | Disk write throughput |
| 4 | Network transmitted throughput |
| 5 | Provisioned capacity for CPU |
| 6 | Use of CPU [MHZ] |
| 7 | Use of CPU [%] |
| 8 | Provisioned Memory capacity [KB] |
| 9 | Memory consumed [KB] |
| 10 | Network received throughput [KB/s] |
| 11 | CPU cores |

TABLE II.  SELECTED FEATURES BY PCGPSONN

| S.No | Feature name |
|---|---|
| 1 | Provisioned capacity for CPU |
| 2 | Use of CPU [MHZ] |
| 3 | Use of CPU [%] |
| 4 | Provisioned Memory capacity [KB] |
| 5 | Memory consumed [KB] |
| 6 | Network received throughput [KB/s] |
| 7 | CPU cores |

TABLE III.  SELECTED ATTRIBUTES BY GPSO

| S.No | Feature name |
|---|---|
| 1 | Disk read throughput |
| 2 | Disk write throughput |
| 3 | Network transmitted throughput |
| 4 | Provisioned capacity for CPU |
| 5 | Use of CPU [MHZ] |
| 6 | Use of CPU  [%] |
| 7 | Memory capacity provisioned [KB] |
| 8 | Memory usage [KB] |
| 9 | Network received throughput [KB/s] |
| 10 | CPU cores |

Table IV displays the GPCNN's chosen attribute list. When using the GPCNN strategy, all attributes' fitness values are first obtained using a genetic algorithm, and they are then given to the PCNN approach for double threshold operation during the attribute selection phase.

TABLE IV.  SELECTED ATTRIBUTES BY GPCNN

| S.No | Feature name |
|---|---|
| 1 | Disk read throughput |
| 2 | Disk write throughput |
| 3 | CPU capacity provisioned |
| 4 | Use of CPU [MHZ] |
| 5 | Use of CPU [%] |
| 6 | Memory capacity provisioned [KB] |
| 7 | Memory usage [KB] |
| 8 | Network received throughput [KB/s] |
| 9 | CPU cores |

## B. Correlation Feature Extraction

The degree to which two or more variables vary together is shown by a statistical metric known as correlation. It describes, in layman's terms, how much one variable changes in reaction to a slight variation in another. In accordance with the direction of the change, it might have positive, negative, or zero values. An independent attribute's strong significance in affecting the output is shown by a high correlation value between it and a dependent attribute. Finding the connection between the dependent and all of the independent variables in a multiple regression setup with numerous parameters is required to produce a more accurate and useful model. It is important to always keep in mind that additional characteristics do not necessarily translate into greater accuracy. Irrelevant characteristics would add unnecessary noise to our model and the accuracy could fall.

The Pearson r correlation method used in this study can be used to identify correlation between two attributes. The Pearson r correlation is the most often used correlation statistic to assess the degree of relationship between two linearly related characteristics. It is possible to determine the Pearson correlation between any two qualities, x, and y using:

$$r_{xy} = \frac{n\Sigma x_i y_i - \Sigma x_i y_i}{\sqrt{n\Sigma x_i^2 - (\Sigma x_i)^2}\sqrt{n\Sigma y_i^2 - (\Sigma y_i)^2}} \qquad (1)$$



Fig. 2. Correlation Feature extracted values using Pearson correlation method.

Pearson correlation coefficient, often denoted as r, that calculate the Pearson correlation coefficients between each selected attributes which are derived from section 3.1 and the target variable (cloud resource demand). Using Pearson correlation coefficient, the extracted correlation feature table is shown in Fig. 2. In the Fig. 2, the Pearson correlation ranges from -1 to 1, where:

- r=1 indicates a perfect positive linear relationship.
- r=−1 indicates a perfect negative linear relationship.
- r=0 indicates no linear relationship.

These extracted correlation feature values are given into the SVM to predict the cloud resource demand.

## C. Cloud Resource Demand Prediction using SVM

The equations are an exception to the prescribed specifications of this template. You will need to determine w Regression analysis and classification are two applications of support vector machines, commonly referred to as support vector networks or SVMs, which are supervised learning models in machine learning. They examine data and identify patterns. Provided a set of training examples that have been labelled as falling into one of two categories, an SVM training process builds a model that places new examples into either group. The model is now a non-probabilistic binary linear classifier as a result of this. The objective of an SVM model is to generate as big of a gap as possible between the instances of the different categories by mapping the examples as points in space. Next, by mapping them into the same region and identifying which side of the gap they fall into, new instances are projected to fit into a particular category. SVMs can perform non-linear classification as well as linear classification by implicitly translating their inputs into feature spaces with many dimensions, a method known as the "kernel trick."

The primary objective of SVM in cloud resource demand prediction is to find a hyperplane that best separates data points belonging to different classes. The hyperplane is chosen to maximize the margin, which is the distance between the hyperplane and the nearest data points from each class, known as support vectors. The hyperplane serves as the decision boundary that separates data points into different classes. The margin in SVM is the distance between the hyperplane and the nearest data points from each class (support vectors). SVM aims to maximize this margin. Support vectors are the data points that lie closest to the decision boundary. They play a crucial role in defining the optimal hyperplane. Train an SVM model in this work, using correlation features as input and cloud resource prediction value as the target variable. Use the trained SVM model to predict cloud resource prediction value for new, unseen data.

The work flow of SVM for cloud resource prediction is shown in Fig. 3.



Fig. 3. SVM algorithm.

## IV. RESULTS

The dataset used for designing and developing the cloud resource demand prediction using SVM comprises information on resource utilization over a period of one month. In total, there were 123 million observed instances involving 1250 machines. The dataset which is used in this work is Real-time dataset named GoCJ. The GoCJ datasets used in our work is the publicly available dataset collected from https://data.mendeley.com/datasets/b7bp6xhrcd/1. The new proposed approach is implemented in Matlab.

### A. Examination Parameters

A wide range of assessment criteria are available to assess the prediction algorithms' performance. This work considers the following elements.

Detection Accuracy is:

$$\frac{True\ Positive(TP)+True\ Negative(TN)}{TP+False\ Positive(FP)+TN+False\ Negative(FN)} \quad (2)$$

Error Rate is:

$$\frac{No\ of\ Datas\ of\ falsely\ labeled\ texts}{Total\ No\ of\ texts} \quad (3)$$

Precision rate:

$$\frac{TP}{TP+FP} \quad (4)$$

Recall rate:

$$\frac{TP}{TP+FP} \quad (5)$$

### B. Performance Analysis

Using the performance measures indicated above, the classifier system's performance is analyzed and contrasted with that of other techniques. The tables and graphs below demonstrate this.

*1) Experiment No #1:* Newly Developed Attribute Selection Approach Evaluation with Accuracy.

We will evaluate the impact of every attribute selection strategy used in the work in this research. Eq. (2) to Eq. (5) are used to evaluate the effectiveness of this cloud resource demand prediction technique. A great attribute selection strategy is preferred to have a high value of Eq. (2). The PCGPSONN accuracy analysis is listed in Table V.

As is evident from Table V, the PCGPSONN has an accuracy of 96.29. The accuracy of attribute selection is shown in Fig. 4.

TABLE V.        NEWLY DEVELOPED ATTRIBUTE SELECTION APPROACH EVALUATION WITH ACCURACY

| GPSO | GPCNN | PCGPSONN |
|---|---|---|
| 93.14 | 94.23 | 96.29 |



Fig. 4.   Accuracy of feature selection strategies.

*2) Experiment No #2:* Newly Developed Attribute Selection Approach Evaluation with Precision Rate.

Precision tells the proportion of instances that the models predicted as positive and were actually positive out of all the instances it predicted as positive. A high precision indicates that when the model predicts a positive outcome, it is likely to be correct. The PCGPSONN Precision rate analysis is listed in Table VI.

TABLE VI.        NEWLY DEVELOPED ATTRIBUTE SELECTION APPROACH EVALUATION WITH PRECISION RATE

| GPSO | GPCNN | PCGPSONN |
|---|---|---|
| 90.21 | 91.35 | 93.10 |

As is evident from Table VI, the PCGPSONN has Precision rate of 93.1, which is better than other approaches. As a result, the PCGPSONN classifier is thought to be the best for choosing attributes. The Precision rate of attribute selection is shown in Fig. 5.



Fig. 5.   Precision rate of attribute selection strategies.

*3) Experiment No #3:* Newly Developed Attribute Selection Approach Evaluation with Recall Rate.

Recall depicts the proportion of actual positive instances that were correctly predicted by the model out of all the actual positive instances. A good attribute selection strategy is preferred to have a high value of Eq. (4). The PCGPSONN Recall rate analysis is listed in Table VII.

TABLE VII.    NEWLY DEVELOPED ATTRIBUTE SELECTION APPROACH EVALUATION WITH RECALL RATE

| GPSO | GPCNN | PCGPSONN |
|------|-------|----------|
| 90.11 | 91.35 | 92.85 |

As is evident from Table VII, the PCGPSONN has Recall rate of 93, which is better than other approaches. The Recall rate of attribute selection is shown in Fig. 6.



Fig. 6.  Recall rate of attribute selection strategies.

As can be seen from the accompanying figure, the PCGPSONN's Recall Rate is better than other methods. Therefore, the PCGPSONN is the best for choosing attributes.

*4) Experiment No #4:* Newly Developed Attribute Selection Approach Evaluation with Error Rate.

Error rate, in the context of machine learning, is a metric that represents the proportion of incorrectly classified instances in a model's predictions. It is the complement of accuracy. The PCGPSONN error rate analysis is listed in Table VIII.

TABLE VIII.    NEWLY DEVELOPED ATTRIBUTE SELECTION APPROACH EVALUATION WITH ERROR RATE

| GPSO | GPCNN | PCGPSONN |
|------|-------|----------|
| 6.86 | 5.77 | 0.03703704 |

As is evident from Table VIII, the PCGPSONN has an error rate of 0.037. The error rate of attribute selection is shown in Fig. 7.

As can be seen from the associated figure, the PCGPSONN's error rate is lower than other methods. Therefore, the PCGPSONN is the best for choosing attributes.



Fig. 7.  Error rate of feature selection strategies.

## V.  CONCLUSION

Demands on cloud resources can be sudden, intense, and volatile. The reactive resource providing approach may result in sluggish or insufficient resource delivery. Thus, to guarantee resource availability, it is imperative to predict resource demands. To process the raw data and produce a fresh and original prediction for Cloud resource demands, machine learning techniques were applied in this study. We were able to develop a more accurate model for cloud resource demand prediction in this study by effectively using a feature selection method based on data mining methodology. PCGPSONN showed to be quite accurate at predicting the demand for Cloud resource. The future direction of this research can be carried out with several machine learning approaches to enhance prediction techniques. Additionally, novel feature selection techniques can be used to develop a deeper comprehension of critical traits and enhance predictions of cloud resource consumption.

## REFERENCES

[1]  M. Ullrich and J. Lässig. 2013. Current Challenges and Approaches for Resource Demand Estimation in the Cloud. 2013 International Conference on Cloud Computing and Big Data, 2013:387-394.

[2]  L. Chen and H. Shen. 2016. Towards Resource-Efficient Cloud Systems: Avoiding Over-Provisioning in Demand-Prediction Based Resource Provisioning. 2016 IEEE International Conference on Big Data, 2016:184-193.

[3]  N. Roy, A. Dubey and A. Gokhale. 2011. Efficient Autoscaling in the Cloud Using Predictive Models for Workload Forecasting. 2011 IEEE 4th International Conference on Cloud Computing, 2011:500- 507.

[4]  Z. Chen, Y. Zhu, Y. Di and S. Feng. 2015. Self-Adaptive Prediction of Cloud Resource Demands Using Ensemble Model and Subtractive Fuzzy Clustering Based Fuzzy Neural Network. Computational Intelligence and Neuroscience, 2015:1-14.

[5]  Y. Jiang, C. Perng, T. Li, et al. 2011. ASAP: A Self-Adaptive Prediction System for Instant Cloud Resource Demand Provisioning. 2011 IEEE 11th International Conference on Data Mining, 2011:1104-1109.

[6]  D. Xu, S. Yang and H. Luo. 2015. Research on Generalized Fuzzy Soft Sets Theory based Combined Model for Demanded Cloud Computing Resource Prediction. Chinese Journal of Management Science, 2015, 23(5):56-64.

[7] M. Verma, G. R. Gangadharan, V. Ravi, et al. 2013. Resource Demand Prediction in Multi-tenant Service Clouds. 2013 IEEE International Conference on Cloud Computing in Emerging Marks, 2013, 28(17):1-8.

[8] S. Imai, T. Chestna and C. A. Varela. 2013. Accurate Resource Prediction for Hybrid IaaS Clouds Using Workload-Tailored Elastic Compute Units. 2013 IEEE/ACM 6th International Conference on Utility and Cloud Computing, 2013:171-178.

[9] H. Mi, H. Wang, G. Yin, D. Shi, Y. Zhou and L. Yuan. 2011. Resource On-Demand Reconfiguration Method for Virtualized Data Centers. Journal of Software, 2011, 22(9):2193-2205.

[10] D. Minarolli and B. Freisleben. 2014. Cross-Correlation Prediction of Resource Demand for Virtual Machine Resource Allocation in Clouds. 2014 Sixth International Conference on Computational Intelligence, Communication Systems and Networks, 2014:119-124.

[11] A. A. Bankole and S. A. Ajila. 2013. Predicting Cloud Resource Provisioning Using Machine Learning Techniques. 2013 26th IEEE Canadian Conference of Electrical and Computer Engineering, 2013:1-4.

[12] K. Rajaram and M. P. Malarvizhi. 2017. Utilization based Prediction Model for Resource Provisioning. 2017 IEEE International Conference on Computer, Communication and Signal Processing, 2017:1-6.

[13] W. Zhang, P. Duan, L. T. Yang, et al. 2017. Resource Requests Prediction in the Cloud Computing Environment with a Deep Belief Network. Software Practice and Experience, 2017, 47(3):473-488.

[14] M. Borkowski, S. Schulte and C. Hochreiner. 2016. Predicting Cloud Resource Utilization. 2016 IEEE/ACM 9th International Conference on Utility and Cloud Computing, 2016:37-42.

[15] Mostafa Bozorgi, S., & Yazdani, S. (2019). IWOA: An improved whale optimization algorithm for optimization problems. Journal of Computational Design and Engineering, 6(3), 243-259.

[16] Manikandan, N., Gobalakrishnan, N., & Pradeep, K. (2022). Bee optimization based random double adaptive whale optimization model for task scheduling in cloud computing environment. Computer Communications, 187, 35-44.

[17] Chen, X.; Cheng, L.; Liu, C.; Liu, Q.; Liu, J.; Mao, Y.; Murphy, J. A WOA-based optimization approach for task scheduling in cloud computing systems. IEEE Syst. J. 2020, 14, 3117–3128.

[18] Madni, S.H.H.; Abd Latiff, M.S.; Abdullahi, M.; Abdulhamid, S.M.; Usman, M.J. Performance comparison of heuristic algorithms for task scheduling in IaaS cloud computing environment. PLoS ONE 2017, 12, e0176321.

[19] Bezdan, T.; Zivkovic, M.; Bacanin, N.; Strumberger, I.; Tuba, E.; Tuba, M. Multi-objective task scheduling in cloud computing environment by hybridized bat algorithm. IFS 2021, 42, 411–423.

[20] Sukhoroslov, O., Nazarenko, A., & Aleksandrov, R. (2019). An experimental study of scheduling algorithms for many-task applications. The Journal of Supercomputing, 75, 7857-7871.

[21] Buang, N.; Hanawi, S.A.; Mohamed, H.; Jenal, R. B-Spline Curve Modelling Based on Nature Inspired Algorithms. APJITM 2016.

[22] Kumar, M., Sharma, S. C., Goel, A., & Singh, S. P. (2019). A comprehensive survey for scheduling techniques in cloud computing. Journal of Network and Computer Applications, 143, 1-33.

[23] Chhabra, A., Huang, K. C., Bacanin, N., & Rashid, T. A. (2022). Optimizing bag-of-tasks scheduling on cloud data centers using hybrid swarm-intelligence meta-heuristic. The Journal of Supercomputing, 1-63.

[24] Ni, L., Sun, X., Li, X., & Zhang, J. (2021). GCWOAS2: multi objective task scheduling strategy based on Gaussian cloud-whale optimization in cloud computing. Computational Intelligence and Neuroscience, 2021, 1-17.

[25] Manikandan, N., Gobalakrishnan, N., & Pradeep, K. (2022). Bee optimization based random double adaptive whale optimization model for task scheduling in cloud computing environment. Computer Communications, 187, 35-44.

[26] Chen, X.; Cheng, L.; Liu, C.; Liu, Q.; Liu, J.; Mao, Y.; Murphy, J. A WOA-based optimization approach for task scheduling in cloud computing systems. IEEE Syst. J. 2020, 14, 3117–3128.

[27] Jia, L., Li, K., & Shi, X. (2021). Cloud computing task scheduling model based on improved whale optimization algorithm. Wireless Communications and Mobile Computing, 2021, 1-13.

[28] Albert, P.; Nanjappan, M. WHOA: Hybrid based task scheduling in cloud computing environment. Wireless. Pers. Communication. 2021, 121, 2327–2345.

[29] Sharma, M.; Garg, R. Energy-aware whale-optimized task scheduler in cloud computing. In Proceedings of the 2017 International Conference on Intelligent Sustainable Systems (ICISS), Palladam, India, 7–8 December 2017; pp. 121–126.

[30] Sreenu, K.; Sreelatha, M. W-Scheduler: Whale optimization for task scheduling in cloud computing. Cluster Computing. 2019, 22, 1087–1098.

[31] Rekha, P.M.; Dakshayini, M. Efficient task allocation approach using genetic algorithm for cloud environment. Cluster Computing. 2019, 22, 1241–1251.

[32] Sun, Y., Li, J., Fu, X., Wang, H., & Li, H. (2020). Application research based on improved genetic algorithm in cloud task scheduling. Journal of Intelligent & Fuzzy Systems, 38(1), 239-246.

[33] Abdullahi, M., & Ngadi, M. A. (2016). Symbiotic organism search optimization based task scheduling in cloud computing environment. Future Generation Computer Systems, 56, 640-650.

[34] Li, G.; Wu, Z. Ant Colony Optimization Task Scheduling Algorithm for SWIM Based on Load Balancing. Future Internet 2019, 11, 90.

[35] Huang, X.; Li, C.; Chen, H.; An, D. Task scheduling in cloud computing using particle swarm optimization with time varying inertia weight strategies. Cluster. Computing. 2020, 23, 1137–1147.

[36] Nabi, S.; Ahmad, M.; Ibrahim, M.; Hamam, H. AdPSO: Adaptive pso-based task scheduling approach for cloud computing. Sensors 2022, 22, 920.

[37] Chen, X.; Long, D. Task scheduling of cloud computing using integrated particle swarm algorithm and ant colony algorithm. Cluster. Computing. 2019, 22, 2761–2769.

[38] Kumar, M., & Sharma, S. C. (2018). PSO-COGENT: Cost and energy efficient scheduling in cloud environment with deadline constraint. Sustainable Computing: Informatics and Systems, 19, 147-164.

[39] Zhou, Z.; Li, F.; Abawajy, J.H.; Gao, C. Improved PSO Algorithm Integrated with Opposition-Based Learning and Tentative Perception in Networked Data Centers. IEEE Access 2020, 8, 55872–55880.

[40] Madni, S. H. H., Latiff, M. S. A., Ali, J., & Abdulhamid, S. I. M. (2019). Multi-objective-oriented cuckoo search optimization-based resource scheduling algorithm for clouds. Arabian Journal for Science and Engineering, 44, 3585-3602.

[41] Madni, S.H.H.; Abd Latiff, M.S.; Abdulhamid, S.M.; Ali, J. Hybrid gradient descent cuckoo search (HGDCS) algorithm for resource scheduling in IaaS cloud computing environment. Cluster Computing. 2019, 22, 301–334.

[42] Chhabra, A.; Singh, G.; Kahlon, K.S. Multi-criteria HPC task scheduling on IaaS cloud infrastructures using meta-heuristics. Cluster. Computing 2021, 24, 885–918.

[43] Natesha, B.V.; Kumar Sharma, N.; Domanal, S.; Reddy Guddeti, R.M. GWOTS: Grey Wolf Optimization Based Task Scheduling at the Green Cloud Data Center. In Proceedings of the 2018 14th International Conference on Semantics, Knowledge and Grids (SKG), Guangzhou, China, 12–14 September 2018; pp. 181–187.

[44] Alzaqebah, A.; Al-Sayyed, R.; Masadeh, R. Task Scheduling based on Modified Grey Wolf Optimizer in Cloud Computing Environment. In Proceedings of the 2nd International Conference on new Trends in Computing Sciences (ICTCS), Amman, Jordan, 9–11 October 2019; pp. 1–6.

[45] Natesan, G.; Chokkalingam, A. Task scheduling in heterogeneous cloud environment using mean grey wolf optimization algorithm. ICT Express 2019, 5, 110–114.

[46] Elaziz, M.A.; Xiong, S.; Jayasena, K.P.N.; Li, L. Task scheduling in cloud computing based on hybrid moth search algorithm and differential evolution. Knowl.-Based Syst. 2019, 169, 39–52.

[47] Srichandan, S., Kumar, T. A., & Bibhudatta, S. (2018). Task scheduling for cloud computing using multi-objective hybrid bacteria foraging algorithm. Future Computing and Informatics Journal, 3(2), 210-230.

[48] Nasr, A.A.; Chronopoulos, A.T.; El-Bahnasawy, N.A.; Attiya, G.; El-Sayed, A. A novel water pressure change optimization technique for solving scheduling problem in cloud computing. Cluster. Computing. 2019, 22, 601–617.

[49] Praveen, S.P.; Rao, K.T.; Janakiramaiah, B. Effective Allocation of Resources and Task Scheduling in Cloud Environment using Social Group Optimization. Arab. J. Sci. Eng. 2018, 43, 4265–4272.

# Spatial Display Model of Oil Painting Art Based on Digital Vision Design

QiongYang[*], Zixuan Yue

School of Architecture Engineering, Xuzhou College of Industrial Technology, Xuzhou, 221000, China

*Abstract*—Oil painting, owing to its unique expressive approach, holds infinite charm in classical artistic creation, yet introduces complexities in terms of manual maintenance. In pursuit of digital spatial visualization of oil painting art, this study employs a stereo matching algorithm and Efficient large-scale stereo matching, focusing on aspects like disparity maps and pixel contrasts. Furthermore, enhancements in the algorithm involve the incorporation of the cross-arms strategy for image registration and the selection of auxiliary point sets to optimize the handling of image features. Results indicate that the proposed model, evaluated on the Middlebury dataset, achieves high accuracy, recall rates, and F1 scores, measuring 97.2%, 95.0%, and 97.5% respectively, surpassing the DecStereo algorithm by 3.4%, 8.2%, and 5.7%. When tested on the Photo2monet oil painting dataset, the proposed model achieves peak signal-to-noise ratio and average structural similarity index values of 16.781 and 0.833 respectively. This suggests that the proposed model excels in digital visual representation of oil paintings, exhibiting higher image precision, stronger stereo matching capabilities, and superior spatial display performance.

*Keywords*—*Oil painting; spatial visualization; Stereo matching; Spatial display; ELAS*

## I. INTRODUCTION

Oil painting is an art form that combines lines and colors, showcasing significant changes in color during creation, primarily reliant on the movement and intensity control of painting tools, employing layered pigment transitions to portray various color combinations [1]. Although colors may change during the creation of an oil painting, it essentially relies on precise control of the painting tools and changes in the layers of pigment to show different color effects. In recent years, with the development of technology, especially advances in digital visual design, the way oil painting art is displayed and created is undergoing a revolutionary change. Three-dimensional visualization methods demonstrate image content from multiple angles through 3D scanning and associated nodal information. Techniques and artistic styles hold pivotal positions in contemporary oil painting creation [2]. However, limitations in traditional oil painting methods and tools lead to susceptibility to damage or color fading, even with strict protective measures. Stereo matching technology, a critical component of stereo vision systems, has garnered increased attention and research among scholars [3]. This method rectifies positional differences between left and right images after epipolar line correction, determining matching cost values based on disparities among image points. As it utilizes partial window regions for computation, it boasts higher computational efficiency and suitability for parallel

processing, especially in scenarios demanding high real-time performance [4]. Yet, due to its reliance on partial image windows, it yields numerous false matches, diminishing accuracy in disparity determination. Three-dimensional solid matching methods eschew cost aggregation, exclusively computing matching costs, disparities, and undergoing post-processing. Overall, employing holistic visual perception methods that consider individual pixel points across the entire image constructs an energy function to optimize disparity estimation for the whole image, thereby enhancing matching accuracy and overall visual effect [5]. Presently, scholars worldwide focus on enhancing image processing accuracy and real-time performance. However, no method adequately balances disparity and matching speed, where high-precision algorithms often sacrifice computational speed. This research focuses on developing a new type of spatial matching model for oil painting art images by applying stereo matching technology and three-dimensional visualization construction methods to enhance the presentation skills and styles of modern oil paintings. Traditional oil painting creation and display are limited by physical media and tools, and are prone to problems such as damage or color loss. In addition traditional methods often fail to provide an interactive and immersive art experience. Therefore, the research is dedicated to overcoming these limitations by utilizing digital technology to accurately capture and analyze the oil brush paths as well as the key frames of the paintings through 3D visualization techniques, so as to preserve and enhance the style and realism in the creation of contemporary oil paintings. In view of this, the stereo matching technique and 3D entity matching method used in this study can efficiently deal with the parallax problem in images and improve the matching accuracy. The use of this technique not only improves the accuracy and real-time of image processing, but also makes the digital presentation of art works more realistic and dynamic. A spatial matching model for oil painting art images based on a fast and effective stereo matching algorithm is developed. The model not only preserves the style and realism of contemporary oil painting, but also enhances the audience's participation and the expressiveness of the artwork. In addition, the research has innovatively transformed the way of displaying and appreciating oil painting art by integrating advanced technologies such as image processing, virtual reality and interaction design, which has brought a long-term impact on the art field. A spatial matching model for oil painting art images based on a fast and effective stereo matching algorithm is studied and developed. The model not only preserves the style and realism in the creation of contemporary oil paintings, but also enhances the audience's participation and the

expressiveness of the artwork. In addition, by integrating advanced technologies such as image processing, virtual reality and interaction design, the research has innovatively transformed the way of displaying and appreciating oil paintings, which has brought a long-term impact on the art field.

The research is mainly divided into four sections: Section II involves a literature review related to digital visual design and art design; Section III focuses on the construction of an oil painting spatial display model based on image stereo matching algorithms; Section III analyzes the performance results and application outcomes of the image stereo matching algorithm; discussion and conclusion is given in Section V and VI respectively.

## II. RELATED WORKS

Numerous scholars have conducted diverse research in the field of digital visual design. Logeshwaran et al. proposed a segmentation-based visual processing algorithm aimed at enhancing resolution and clarity to some extent through multi-visual enhanced pixels [6]. Parker and others studied the impact of digital technology on artwork design, illustrating how new technologies positively and negatively affect work resources based on various factors [7]. Peng and colleagues considered visibility analysis as an important research area in visual landscape research and developed a new computational algorithm for complex environments that can analyze views from multiple angles [8]. Johnson et al. constructed a pseudo-image-based color mapping table, optimized 3D scanning meshes for data visualization, and synthesized textures from pseudo-images. Additionally, there is an interactive rendering engine with custom algorithms and interfaces that can showcase various new visual styles to depict points, lines, surfaces, and volumetric data [9]. Ye et al. proposed a feasible method to quantitatively measure perception-based street visual quality and developed a Java-based program to automatically collect experienced urban designers' preferences for representative sample images. The results indicated that the evaluation model has satisfactory accuracy and provides insights into the perception-based visual quality and detailed mapping of various key elements in streets, offering accurate design guidance to aid more effective street renovations [10]. Chen et al. achieved a paradigm shift in perceptive environments by capturing local pixel-level intensity changes and generating asynchronous event streams. From standard computer vision to event-based neural morphological vision, advanced technologies in autonomous vehicle visual sensing systems have been developed [11]. Ren et al. systematically discussed the history and current status of optical lighting, image acquisition, image processing, and image analysis in the field of visual detection. The latest advances in machine vision-based industrial defect detection were introduced, emphasizing the increasing importance of deep learning in the further development of visual detection [12].

In the realm of digital research in art, particularly in the digitization of oil paintings, Castellano G and others have outlined some of the most relevant deep learning methods for pattern extraction and recognition in visual arts, especially in painting and sketching. The ongoing advancements in deep learning and computer vision, coupled with the ever-expanding collections of digitally visualized artworks, present new opportunities for computer science researchers [13]. Sandoval and colleagues introduced a novel two-stage image classification method aimed at enhancing the accuracy of style classification. Experiments conducted on three standard art classification datasets indicated a significant improvement over existing baseline techniques [14]. Scholar Mao W conducted research and analysis on oil painting art education videos using machine learning combined with virtual reality computing. Through the application of virtual reality technology in teaching practices, the study analyzed the effectiveness of teaching methods, allowing students to immerse themselves in art appreciation activities, embrace multiculturalism, experience learning, and enhance aesthetic qualities [15]. Mills and others investigated users' creative digital designs, including a popular 3D virtual painting program. The analysis focused on how students convey the same story in written, oral, and virtual painting modes, tracking key themes in students' virtual experiences [16]. Scholar Kent examined the works of McNeill and Hamill, emphasizing the importance of conducting abstract experiments in virtual reality as a creative medium [17]. Scholar van der Veen M approached the intersection of two visual modes from an augmented reality perspective. In this approach, the real environment is sensed by machines, magnified, scanned, located, and linked to a 3D model [18]. Scholar Doyle provided a historical overview of virtual reality, analyzed major works in the first wave, and discussed the application of virtual reality in contemporary practices through the concept of emotional engagement [19]. Ren Shihong et al. proposed a web-based VPL, JSPatcher, in order to address the problem of multimedia presentation or content generation.The tool allows users to build audio graphs using the Web Audio API and to graphically design and run digital signal processor algorithms using domain-specific audio processing languages. Experimental results show that the tool can be used with other JavaScript language built-ins, Web APIs, etc. to create interactive programs and artworks that can be shared online [20]. Nicole Johnson-Glauch et al. in order to explore how students can utilize representational features, proposed a synthesis of the results of two previous studies Method. The results suggest that visual representations have an impact on students' ability to access and use domain knowledge. Students may confuse concepts represented by similar features and not use concepts without salient features, and statics students are more coordinated in switching representations. These findings provide a generic domain pathway for redesigning the notation and representation of engineering concepts and suggest future research directions [21].

In summary, scholars have primarily focused their attention on environmental perception or programming methods in the digital visualization of images. There is a limited application of these approaches in the direction of oil painting. The integration of artworks and computer vision is often approached through theoretical analysis with a lesser emphasis on incorporating machine vision algorithms. In light of this, the current study aims to design a spatial display

model for oil painting art based on an improved image stereo matching algorithm.

### III. Oil Painting Art Spatial Display Model Based on Improved Stereo Matching Algorithm

To achieve the three-dimensional spatial display of oil painting art, this study, from the perspective of image stereo matching, constructs a method of digital visualization for oil paintings based on a fast and efficient stereo matching algorithm. The goal is to realize the modern intelligent preservation of traditional hand-painted oil paintings.

#### A. Stereo Matching Algorithm for Oil Painting Images Based on PatchMatchStereo-ELAS

PatchMatch Stereo This is an image matching algorithm used in stereo vision. The core idea of PatchMatch is to find the best match between images using random search. This method finds correspondences between images by quickly and randomly guessing the matching points and iteratively improving them. PatchMatch Stereo is used to compute parallax maps between images captured by two cameras, which is very important for reconstructing the three-dimensional structure of a scene. Efficient Large-Scale Stereo Matching (ELAS ) algorithm is an efficient stereo matching method mainly used to compute the depth information of large-scale scenes. It employs a dense method that utilizes geometric constraints between pixels to improve the matching process, thus improving the computational efficiency and matching accuracy. Drawing inspiration from the PatchMatch approach, it employs various propagation techniques such as random initialization, spatial propagation, viewpoint propagation, temporal propagation, and surface optimization to obtain sub-pixel disparity values [22]. However, this method requires the initialization of unit normal vectors and disparity values for each pixel in an oil painting image. A challenge arises as the chosen initial values often deviate significantly from the actual data, making it difficult to find a surface that best fits the real data during the iterative algorithm, resulting in prolonged computation times. In this study, an efficient large-scale stereo matching algorithm (ELAS) is proposed to acquire initial information for oil painting images and 3D labels, thus circumventing the issue of arbitrary initial values in traditional three-dimensional image processing. Fig. 1 illustrates the process of the Patch Match Stereo-ELAS algorithm.

The PatchMatchStereo algorithm initially initializes the disparity map and surface normals for each pixel in oil painting images and then iteratively propagates the disparity map and surface normals. This process incurs a substantial computational cost. On the other hand, ELAS divides the problem into multiple disparity surfaces to obtain robust matching support points. While it enables rapid computation of the disparity map, its performance on low-dimensional measurement images is often inferior to high-resolution images. To address this, the proposed approach first utilizes the ELAS method to obtain robust support points and disparity values. Subsequently, the original disparity values and the constructed parameters of the disparity surface are input into the PatchMatchStereo algorithm. This joint algorithm utilizes support points as vertices within the smaller texture range of

oil painting images, enhancing the disparity effect. Eq. (1) represents the calculation of pixel disparity values.

$$d_p = a_{f_p} p_u + b_{f_p} p_v + c_{f_p} \tag{1}$$

In Eq. (1), $d_p$ represents the disparity value, $a_{f_p}$, $b_{f_p}$, $c_{f_p}$ represent disparity parameters, and $p$ represents a pixel. With the aid of equation 1, the problem of estimating disparity can be transformed into the estimation of the disparity plane, i.e., determining how to satisfy the optimization conditions of the disparity surface for each pixel. Eq. (2) illustrates the relationship between the three parameters of the disparity plane and the plane's normal vector.

$$\begin{cases} a_{f_p} = -n_x / n_z \\ b_{f_p} = -n_y / n_z \\ c_{f_p} = -\left(n_x * p_u + n_y * p_v + n_z * d_p\right)/ n_z \end{cases} \tag{2}$$

In Eq. (2), $n_x$, $n_y$, and $n_z$ represent the coordinate values in the $xyz$ directions. Equation 3 is the expression for computing the matching cost.

$$m(p,f) = \sum_{q \in W_p} w(p,q)\beta\left(q, q - \left(a_f q_u + b_f q_v + c_f\right)\right) \tag{3}$$

In Eq. (3), $m(p,f)$ represents the cost value, $f$ represents the disparity plane, $W_p$ represents the specified square matching window, and $w(p,q)$ represents the adaptive weight. Eq. (4) is the expression for calculating the adaptive weight.

$$w(p,q) = e^{-\left\|I_p - I_q\right\|/\gamma} \tag{4}$$



Fig. 1. Process diagram of patchmatchstereo ELAS algorithm.

In Eq. (4), $\gamma$ represents a custom parameter, and $\left\| I_p - I_q \right\|$ represents the L1 distance of the RGB colors of pixels $p$ and $q$. If the color difference between two adjacent pixels is large, then these two points are less likely to lie on the same plane, hence $w(p,q)$ is smaller. The Census transformation compares the grayscale values of the central pixel in a given window with the grayscale values of adjacent pixels. The size relationship is set to 0 and 1, expressing the grayscale information of the central pixel as a binary bit sequence of 0s and 1s. By comparing the bits, the Hamming distance is calculated. Fig. 2 illustrates the Census transformation.



Fig. 2. Schematic diagram of census transformation.

Census transformation effectively considers the differences between points in the window and the central point, reducing the adverse effects of changes in lighting conditions. It is robust in regions of oil painting images with low noise and weak texture. In this process, the study will use both color and grayscale information in the image, combining them with Census transformation for image matching. Eq. (5) represents the improved matching cost expression.

$$\beta(q,q') = \frac{(1-\alpha)}{1+\varepsilon}\left(\left\|G_q - G_{q'}\right\| + \varepsilon * \left\|Cen_q - Cen_{q'}\right\|\right) + \alpha \min\left(\left\|\nabla I_q - \nabla I_{q'}\right\|, \tau_{grad}\right) \quad (5)$$

In Eq. (5), $\left\|G_q - G_{q'}\right\|$ represents the absolute difference in grayscale values between pixels $q$ and $q'$, $\left\|Cen_q - Cen_{q'}\right\|$ is the Hamming distance of their Census values, $\varepsilon$ is a custom scalar parameter for adjusting the ratio of feature grayscale size to Census transformation, and $\left\|\nabla I_q - \nabla I_{q'}\right\|$ represents the absolute difference in gradients between $q$ and $q'$. $\alpha$ is a custom scalar parameter for balancing the color and gradient proportions of each pixel. Fig. 3 shows the spatial propagation schematic of the PMS-ELAS algorithm.

The PMS-ELAS algorithm utilizes the prior disparity information obtained from ELAS to enhance the positioning accuracy in weak texture areas and the local optimization ability in weak or non-textured regions. It iteratively goes through two phases: spatial expansion and view propagation,

involving odd and even iterations in both upward and downward directions. This approach can be applied to multiple images and reduces the disparity regions in PMS, thereby improving computational speed and potentially eliminating the need for plane optimization.



Fig. 3. Schematic diagram of spatial propagation of PMS-ELAS algorithm.

### B. ELAS Algorithm Improved by Cross-Arm

The ELAS algorithm utilizes pixel counts with high confidence levels to form a triangular disparity plane, followed by rapid interpolation on the triangular surface. This method not only reduces computational speed but also swiftly calculates disparities in high-resolution images. However, mismatches in support points can introduce errors in constructing the disparity search distance, increasing the disparity values between points on the disparity surface. Consequently, this diminishes the disparity after interpolation and maximum a posteriori probability estimation on the same plane. Moreover, since the feature points extracted from oil paintings are often distributed in texture-rich and rapidly changing areas, ELAS-derived support points are mostly located at borders and areas with drastic changes. Robust support points are scarce for weakly textured or non-textured regions. To address these issues, research is proposed to enhance the accuracy and efficiency of support point localization and improve the accuracy of disparity estimation using an ELAS algorithm based on cross-arm enhancement. Fig. 4 illustrates the schematic diagram of the improved ELAS framework.



Fig. 4. Schematic diagram of improving the ELAS framework.

Considering the relationship between a single oil painting image and hierarchical images, the disparity calculation for each hierarchical image is based on support point selection and matching costs. Cross-arm structures are employed for image registration and auxiliary point set selection. By matching hierarchical images, confidence values for each level

are obtained, enabling layered processing of images. Sampling is primarily conducted in a pyramid-like fashion, and based on confidence values, support point sets for high-resolution images are updated, simultaneously reducing the search range for each hierarchical image. Using a 3D Sobel operator on both left and right images, the grayscale gradients in the X and Y directions are employed for grayscale estimation in the left and right images. Feature descriptors for the images are constructed by selecting appropriate gradient values in the vicinity of each pixel among 32 points and computing corresponding matching costs. Let $(u_n, v_n)$ be the pixel coordinates, $d(u_n, v_n)$ D be the disparity of the pixel, and $f(u_n, v_n)$ be the 32-dimensional feature vector in the region, including gradients in the horizontal and vertical directions for each pixel. Fig. 5 illustrates the selection method for auxiliary points.



Fig. 5. Selection method of auxiliary points.

The local matching algorithm utilizes a larger window width, resulting in higher image quality and reduced noise. However, it tends to lose texture details and boundary features in the image while also leading to increased time complexity and decreased efficiency. On the other hand, a smaller window can capture more texture details and edge information, providing faster operation and effectively extracting image features to enhance algorithm execution efficiency. By jointly processing multiple feature vectors of multiple points, referred to as auxiliary points, and selecting these auxiliary points strategically, multiple small windows in different regions can be merged into larger windows, leveraging their respective advantages. This improves the overall accuracy of image registration while minimizing the image size for efficient computation, emphasizing local features to enhance computational efficiency. The proposed cross-shaped structure extends the arms to increase the selection window for auxiliary points, reducing dual inconsistencies. Establishing a window involves determining the arm lengths in the upward, downward, leftward, and rightward directions from the central pixel, with the stretch length determined based on the pixel differences. Stretching occurs in a specific direction until a significant color difference is encountered, and there is a maximum arm length limit. The cross arms maximize the perception range of the recognition window and efficiently detect boundaries. Eq. (6) represents the mathematical expression for confidence calculation.

$$cf^s_{(u,v)} = 1 - \frac{\min_d E(u,v,d)}{\min_{d \neq D(u,z)} E(u,v,d)} \qquad (6)$$

In Eq. (6), $E(u,v,d)$ represents the matching cost value of pixel $(u,v)$ at disparity $d$. $D(u,v)$ represents the disparity when the minimum cost value is found in the disparity search area. Confidence is constrained to the range [0,1]. The ratio between the minimum and second-minimum cost values reflects confidence. Fig. 6 illustrates the schematic diagram of image pyramid upsampling.



Fig. 6. Schematic diagram of image pyramid upsampling.

In local stereo matching, effective local information around the central pixel improves matching accuracy. For smaller texture areas, increasing the cross configuration expands alignment coverage, captures more image features, and enhances feature characterization accuracy. However, using a larger window width increases computational complexity and noise. Therefore, a method combining image pyramids and confidence is proposed for multi-scale fusion of support point sets. Downsampling two stereo image pairs yields a higher spatial resolution image pair, from which support points and corresponding disparity and confidence at different levels are obtained. If the disparity values of pixels satisfy left-right consistency conditions, their confidence values increase by a constant $\triangle cf$. In high-resolution images, a deep downsample of the disparity and confidence maps is performed to obtain the disparity map and confidence of the high-resolution image. Upsampling is then used to obtain a higher confidence support point set when searching for support points.

## IV. ANALYSIS OF THE UTILITY OF THE OIL PAINTING ARTISTIC SPACE DISPLAY MODEL BASED ON IMPROVED STEREO VISION MATCHING ALGORITHM

To validate the performance of the stereo matching based on the improved ELAS algorithm on image datasets and evaluate its capability in stereo matching spatial imaging applications, this study selected the Middlebury dataset, KITTI dataset, and Photo2monet oil painting style dataset as experimental samples. Diverse image quality evaluation metrics were employed for data analysis, resulting in comparative results for different image recognition algorithms.

### A. Performance Analysis of the Improved Stereo Vision Matching Algorithm

The study utilized the Middlebury dataset and KITTI

dataset as experimental samples for the visual matching algorithm. The Middlebury dataset comprises 33 pairs of static stereo images in indoor settings, with a maximum resolution of up to 6 million pixels and a maximum disparity range from 200 to 800. The KITTI dataset is used for stereo matching and includes 389 pairs of grayscale stereo images, consisting of 194 training data pairs and 195 testing data pairs. The compared algorithms include the ELAS algorithm, DecStereo algorithm, and Semi-global Matching (SGM). Evaluation metrics include accuracy, recall, and F1 score. Fig. 7 shows the accuracy results.

From Fig. 7(a), it is evident that the proposed model achieved the highest accuracy on the Middlebury dataset, reaching 97.2%. In comparison, the accuracy values of the ELAS algorithm, DecStereo algorithm, and SGM algorithm are all below 96%, with values of 90.0%, 93.8%, and 90.2%, respectively. Therefore, the proposed model improved accuracy on this dataset by 7.2%, 3.4%, and 7.0%. From Fig. 7(a), On the KITTI dataset, the proposed model achieved the highest accuracy of 95.3%. In contrast, the accuracy values of the ELAS algorithm, DecStereo algorithm, and SGM algorithm are all below 92%, with values of 87.8%, 88.5%,

and 91.6%, respectively. Consequently, the proposed model improved accuracy on this dataset by 7.5%, 6.9%, and 3.7%. Fig. 8 presents the recall results.

From Fig. 8, it can be observed that on the Middlebury dataset, the recall curve of the proposed model shows a progressively increasing trend, reaching a convergence value of 95.0% at 50 iterations. In contrast, the recall curves of the ELAS algorithm, DecStereo algorithm, and SGM algorithm exhibit significant fluctuations with a decreasing trend during iterations, and their final convergence values are 85.2%, 86.8%, and 88.6%, respectively. Therefore, the proposed model outperforms them by 9.8%, 8.2%, and 6.4%, respectively. On the KITTI dataset, the recall curve of the proposed model exhibits a stepwise linear increasing trend, reaching a convergence value of 96.8% at 50 iterations. Comparatively, the final convergence values of the ELAS algorithm and SGM algorithm are 87.5% and 90.3%, respectively. The DecStereo algorithm did not achieve a convergence value. Consequently, the proposed model surpasses them by 9.3% and 6.5%. Fig. 9 illustrates the F1 score results.



(a) Middlebury dataset      (b) KITTI dataset

Fig. 7. Accuracy result chart.



(a) Middlebury dataset      (b) KITTI dataset

Fig. 8. Recall rate result chart.

|                      |                  |
| :------------------: | :--------------: |
| (a) Middlebury dataset | (b) KITTI dataset |

Fig. 9.   F1 value result graph.

For the Middlebury dataset, the F1 score curves of the proposed model and the comparative models show fluctuating yet gradually increasing trends. At 50 iterations, the convergence values of the proposed model, ELAS algorithm, DecStereo algorithm, and SGM algorithm are 97.5%, 95.0%, 91.8%, and 91.8%, respectively. Hence, the proposed model outperforms them by 2.5%, 5.7%, and 5.7%. Concerning the KITTI dataset, only the proposed model, DecStereo algorithm, and SGM algorithm achieved convergence values, which are 94.2%, 89.8%, and 89.6%, respectively. Therefore, the proposed algorithm surpasses them by 4.4% and 4.6%, while the ELAS algorithm did not achieve a convergence value.

### B. Analysis of the Application Effect of a Spatial Matching Display Model for Oil Painting Art

The study employed the Photo2monet oil painting style dataset for experimentation, with a training-to-testing set ratio of 3:1. All dataset images were RGB color images of size 256×256. The evaluation metrics chosen for the study included the mismatch rate, average absolute error, Peak Signal to Noise Ratio (PSNR), and Structural Similarity (SSIM). ELAS algorithm, DecStereo algorithm, and SGM algorithm were selected as comparative algorithms. Fig. 10 illustrates the mismatch rate results.

From Fig. 10, it is evident that, in both the training and testing sets, the proposed model consistently achieved the lowest mismatch rates, with curves showing variations around a horizontal value of 2. The lowest and highest mismatch rates for the training set were 1.0% and 2.5%, respectively. For the testing set, the corresponding values were 0.9% and 3.6%. In contrast, the ELAS and SGM algorithms exhibited mismatch rates above 5% in both training and testing phases, while the DecStereo algorithm fluctuated around 3.5%. The proposed algorithm reduced the fluctuation level by 1.4%. Fig. 11 presents the average absolute error results.

From Fig. 11(a) shows that during the training phase, the proposed model exhibited a smooth descent followed by horizontal convergence, starting to converge at 20 iterations, with a final average absolute error convergence value of 1.78%. In contrast, the average absolute error curves for ELAS, DecStereo, and SGM algorithms in the training dataset displayed significant fluctuations before 25 iterations, with final convergence values of 2.40%, 4.35%, and 2.80%, respectively. Hence, the proposed model reduced the error by 0.62%, 2.57%, and 1.02%. In the testing phase, From Fig. 11(a), the proposed model achieved a convergence value of 1.68%, while the DecStereo and SGM algorithms had corresponding values of 3.96% and 4.62%, indicating reductions of 2.28% and 2.94%, respectively. Table I summarizes the image quality assessment results.



|                 |               |
| :-------------: | :-----------: |
| (a) Training set | (b) Testing set |

Fig. 10. Mismatch rate.

(a) Training set



(b) Testing set

Fig. 11. Mean absolute error.

TABLE I.        IMAGE QUALITY EVALUATION RESULTS

| Oil Painting Dataset | Evaluating indicator | ELAS | DecStereo | SGM | Model in this article |
|---|---|---|---|---|---|
| Training set | PSNR | 11.214 | 14.65 | 16.781 | 18.429 |
| | SSIM | 0.346 | 0.801 | 0.726 | 0.888 |
| Testing set | PSNR | 9.464 | 15.564 | 12.927 | 16.525 |
| | SSIM | 0.301 | 0.833 | 0.622 | 0.846 |

From Table I, it can be observed that the model proposed by the research institute achieved the highest values in image quality evaluation metrics, both in the training and testing sets. Specifically, the PSNR (Peak Signal-to-Noise Ratio) and SSIM (Structural Similarity Index) metrics for the training set were 18.429 and 0.888, respectively, while for the testing set, these metrics were 16.525 and 0.846. In comparison, the highest values for the PSNR and SSIM metrics of the contrast models ELAS, DecStereo, and the SGM algorithm were 16.781 and 0.833, indicating an improvement of 1.648 and 0.042 by the proposed model.

## V.  DISCUSSION

In the field of contemporary art, oil painting, as an ancient and classic art form, has been widely appreciated and studied. However, with the rapid development of digital technology, traditional oil painting art is facing new challenges and opportunities. When comparing the accuracy of the system model, it is found that the accuracy of the algorithmic model used in the research can reach the highest value in the dataset Middlebury, in which for the accuracy of other algorithmic models the accuracy of the research using the algorithmic model is higher, and it can get a better ability to test the data, which indicates that the research using the model of the current construction can achieve better data processing ability on the dataset. At the same time on the KITTI dataset, the accuracy rate of the research use method is also the highest value, which at the same time shows that the research use method in the two datasets to get the best accuracy rate, the research use method for the different datasets have less impact on the accuracy rate of all can get a relatively good accuracy rate. In comparing the recall of different algorithmic models,

the recall curve of the research proposed model shows an increasing trend, while the recall of the research used model is higher. This indicates that the current research use algorithmic model is able to maintain a relatively smooth curve change in the recall test, and has a better data test correctness rate compared to other algorithmic models. When comparing the F1 values of different algorithmic models, the F1 values of the research use model can reach 97.5% and 94.2%, comparing with the other algorithmic models higher F1 values, which indicates that the research use model performs better indicators in different data sets, and it can achieve a better processing effect for the processing of oil painting image art. When comparing the change of the false matching rate of the test set and training set of different algorithm models, the research and use method is at the lowest value of the false matching rate index, and the curves show a straight line around the level value of 2. Meanwhile, the change of the false matching rate of the research and use model is higher than that of the other algorithm models, which is due to the fact that the research and use method can achieve better stability in the training and testing. When comparing the convergence of the algorithmic models, the convergence curve of the research use model shows a smooth decline and then a horizontal convergence trend, and starts to converge when the number of iterations is 20, and the final average absolute error convergence value is 1.78%, then it will show that the research use method has better convergence, and it is easier for convergence to take place, which has a better performance for the testing and analyzing of the data. In comparing the quality assessment of the research use method, the research use method is able to achieve better indicator values, which indicates that the research use method is able to get better data

testing results either in the training or testing set.

In summary, the research using method can perform well in the comparison of different data indexes, which indicates that the data model constructed by the current research can complete and achieve better data analysis and testing in the optimization of image data processing, and the method has better performance, which has a good guiding value for the subsequent research in this direction.

## VI. CONCLUSION

The enduring artistic value of oil paintings has attracted attention in the context of integrating it with the trends of digital modernization. To achieve virtual creation of oil paintings in a three-dimensional visual presentation, this study, The research successfully integrates the artistic value of oil painting with digital modernization technology, and realizes the virtual presentation of oil painting art in three-dimensional visual space through the improved stereo matching model of ELAS algorithm. Meanwhile, the algorithm model algorithm is improved from the level of parallax map and pixel contrast in order to optimize the processing of image features, especially in the cross-crossing arm strategy image alignment and the selection of auxiliary point sets are improved and innovated. This research adopts the stereo matching technology and three-dimensional entity matching method, which not only improves the accuracy and real-time of image processing, but also makes the digital presentation of artworks more realistic and dynamic. Meanwhile, the research also developed a fast and effective stereo matching algorithm for the spatial matching model of oil painting art images, which retains the style and realism in the creation of oil paintings, and enhances the audience's participation and the expressive power of the artwork. The results indicate that, for the KITTI dataset, the proposed model achieved accuracy, recall, and F1 values of 95.3%, 96.8%, and 94.2%, respectively, surpassing the SGM algorithm by 3.7%, 6.5%, and 4.6%. For the Photo2monet oil painting style dataset, the proposed model exhibited the lowest error matching rate, with fluctuations around 2% during both training and testing phases, while the error matching rate for the DecStereo algorithm fluctuated around 3.5%. In terms of the mean absolute error metric, the proposed model achieved a convergence value of 1.68% in the final testing phase, representing a reduction of 2.28% and 2.94% compared to the DecStereo and SGM algorithms, respectively. Therefore, these findings suggest that the proposed spatial stereo matching model has a distinct advantage in the digital presentation of oil painting art, preserving the realism and artistic style of oil paintings. This study has achieved a lot of results, but there are some problems with the use of data in the study, such as the use of a small dataset, for this reason more and larger datasets will be included in subsequent studies to be analyzed.

## REFERENCES

[1] I. Bonaduce, C. Duce, A. Lluveras-Tenorio. Conservation issues of modern oil paintings: a molecular model on paint curing. Accounts of chemical research, 2019, 52(12): 3397-3406.

[2] G. Yang. The imagery and abstraction trend of Chinese contemporary oil painting. Linguistics and Culture Review, 2021, 5(2): 454-471.

[3] X. Zhou, D. In, X. Chen, H. M. Bruhn, X. Liu, Y. Yang. Spectral 3D reconstruction of impressionist oil paintings based on macroscopic OCT imaging. Applied optics, 2020, 59(15): 4733-4738.

[4] Z. Chen, J. Li. Application of Multimedia Data Feature Extraction Technology in Teaching Classical Oil Painting. International Journal of Web-Based Learning and Teaching Technologies (IJWLTT), 2023, 18(2): 1-17.

[5] X. Tie, Y. Goh. Evaluating the Effects of Mastery of Techniques (MT), Painting Materials (PM), Choice of Subject Matter (SM), Teaching Methods (TM) on Teaching Effect of Oil Painting (TE) Using PLS-SEM Approach. International Journal of Education & Technology, 2023, 1(1): 51-63.

[6] J. Logeshwaran, M. Ramkumar, T. Kiruthiga. SVPA-the segmentation based visual processing algorithm (SVPA) for illustration enhancements in digital video processing (DVP). ICTACT Journal on Image and Video Processing, 2022, 12(3): 2669-2673.

[7] S. K. Parker, G. Grote. Automation, algorithms, and beyond: Why work design matters more than ever in a digital world. Applied Psychology, 2022, 71(4): 1171-1204.

[8] Y. Peng, S. Nijhuis. A GIS-based algorithm for visual exposure computation: the west lake in Hangzhou (China) as example. Journal of Digital Landscape Architecture, 2021, 6(1): 424-435.

[9] S. Johnson, F. Samsel, G. Abram, D. Olson, J. Solis A. Artifact-based rendering: harnessing natural and traditional visual media for more expressive and engaging 3D visualizations. IEEE transactions on visualization and computer graphics, 2019, 26(1): 492-502.

[10] Y. Ye, W. Zeng, Q. Shen. The visual quality of streets: A human-centred continuous measurement based on machine learning algorithms and street view images. Environment and Planning B: Urban Analytics and City Science, 2019, 46(8): 1439-1457.

[11] G. Chen, H. Cao, J. Conradt, H. J. Tang, F. Rohrbein. Event-based neuromorphic vision for autonomous driving: A paradigm shift for bio-inspired visual sensing and perception. IEEE Signal Processing Magazine, 2020, 37(4): 34-49.

[12] Z. Ren, F. Fang, N. Yan, Y. Wu. State of the art in defect detection based on machine vision. International Journal of Precision Engineering and Manufacturing-Green Technology, 2022, 9(2): 661-691.

[13] G. Castellano, G. Vessio. Deep learning approaches to pattern extraction and recognition in paintings and drawings: An overview. Neural Computing and Applications, 2021, 33(19): 12263-12282.

[14] C. Sandoval, E. Pirogova, M. Lech. Two-stage deep learning approach to the classification of fine-art paintings. IEEE Access, 2019, 7(1): 41770-41781.

[15] W. Mao. Video analysis of intelligent teaching based on machine learning and virtual reality technology. Neural Computing and Applications, 2022, 34(9): 6603-6614.

[16] K. A. Mills, A. Brown. Immersive virtual reality (VR) for digital media making: transmediation is key. Learning, Media and Technology, 2022, 47(2): 179-200.

[17] C. Kent. Beyond Representation in Virtual Reality: The Abstract Art of Jane LaFarge Hamill and Kevin Mack. Leonardo, 2022, 55(3): 240-245.

[18] M. van der Veen. Crossroads of seeing: about layers in painting and superimposition in Augmented Reality. AI & SOCIETY, 2021, 36(4): 1189-1200.

[19] D. Doyle. The two waves of virtual reality in artistic practice. Virtual Creativity, 2021, 11(2): 189-206.

[20] S. Ren, L. Pottier, Buffa M, Yu Yang. JSPatcher, a Visual Programming Environment for Building High-Performance Web Audio Applications. Journal of the Audio Engineering Society, 2022, 70(11):938-950.

[21] Nicole Johnson-lauch, D. S. Choi, G. Herman. How engineering students use domain knowledge when problem : olving using different visual representations. Journal of Engineering Education, 2020, 109(5):443-469.

[22] S. Choudhuri, S. Adeniye, A. Sen. Distribution Alignment Using Complement Entropy Objective and Adaptive Consensus-Based Label Refinement For Partial Domain Adaptation. Artificial Intelligence and Applications. 2023, 1(1): 43-51.

# Research on Neural Network-based Automatic Music Multi-Instrument Classification Approach

Ribin Guo

Shanxi College of Applied Science and Technology, Taiyuan, Shanxi 030062, China

*Abstract*—The automatic classification of multi-instruments plays a crucial role in providing services for music retrieval and recommendation. This paper focuses on automatic multi-instrument classification. Firstly, instrument features were analyzed, and Mel-frequency cepstral coefficient (MFCC) and perceptual linear predictive coefficient (PLPC) were extracted from instrument signals. Features were selected using the entropy weight method. The optimal initial weight threshold of a back-propagation neural network (BPNN) was obtained by utilizing the sparrow search algorithm (SSA), achieving a SSA-BPNN classifier. Experiments were conducted using the IRMAS dataset. The results demonstrated that the combination of MFCC and PLPC selected through the entropy weight method achieved the best performance in automatic multi-instrument classification. The method yielded high P value, recall rate, and F1 value, 0.72, 0.71, and 0.71, respectively. Moreover, it outperformed other algorithms such as support vector machine and XGBoost. These results confirm the reliability of the automatic multi-instrument classification method proposed in this paper, making it suitable for practical applications.

*Keywords*—*Neural network; musical instrument; automatic classification; auditory feature; sparrow search algorithm*

## I. INTRODUCTION

Music serves as a medium for conveying emotions [1]. With the continuous progress and popularization of computer technology, an increasing amount of music information circulates and spreads on the internet, offering users a more convenient means to enjoy music. In order to enhance the user experience further, music information retrieval (MIR) has gradually emerged as a crucial area of research [2]. Through MIR, users can efficiently discover their preferred music. MIR research includes the identification and classification of musical instruments, genres, and styles [3]. Automatic classification of musical instruments refers to the use of intelligent algorithms to automatically classify different musical instruments through processing their signals. In multi-instrument automatic classification, different instrument signals can easily interfere with each other, resulting in a decrease in classification effectiveness.

Neural networks are a type of machine learning that possess strong self-learning capabilities, enabling them to derive useful conclusions from a series of complex and seemingly unrelated data. They have wide applications in fields such as speech recognition, image processing, and automation control, where they demonstrate certain advantages in handling complex data. Therefore, in order to further improve the accuracy of multi-instrument automatic classification, this study focuses on researching neural networks.

Based on the analysis of musical instrument signals, this article selected back-propagation neural network (BPNN) as the main algorithm. In order to further improve its performance, the sparrow search algorithm (SSA) was used to optimize BPNN. Finally, the performance of the proposed method was verified on the IRMAS dataset.

The main contribution of this article lies in providing a more accurate method for automatic classification of multiple musical instruments, which also offers some new references for music signal classification and even speech signal classification. This is conducive to promoting the further application of neural network algorithms in the field of music research.

## II. LITERATURE REVIEW

With the continuous development of machine learning and other technologies, an increasing number of algorithms have been applied in MIR research.

In terms of music popularity prediction, Martin-Gutierrez et al. [4] introduced a deep learning architecture called HitMusicNet for predicting the popularity of music recordings. The experimental results demonstrated its superior predictive capabilities compared to existing techniques. Voetter et al. [5] presented two novel datasets to predict song popularity based on the data from AcousticBrainz, Billboard Hot 100, the Million Song dataset, and last.fm. They verified the usability of the designed dataset by performing experiments on different models.

In terms of singer recognition, Rajesh and Nalini [6] conducted experiments to validate the effectiveness of their approach in singer recognition, which involved the integration of MFCC with chroma DCT-reduced pitch features. Kooshan et al. [7] developed a singer recognition system by integrating deep learning and feedforward neural networks. They utilized long short-term memory (LSTM) to identify vocal frames within music segments, followed by the application of another LSTM for singer identification. Experimental validation confirmed the superior performance of this approach compared to existing methods.

In terms of music recommendation, Feng et al. [8] proposed to model and combine melody, chord, and rhythm features and utilized a multilayer perceptron for music recommendation. The experimental results indicated a 3.52% improvement over the best baseline. Elbir et al. [9] developed an innovative deep neural network that utilizes the acoustic attributes of music to classify genres and provide music recommendations.

In the domain of instrument recognition and classification, Chatterjee et al. [10] proposed employing a convolutional siamese networks along with its residual variation to identify instruments based on scalograms derived from audio recordings. They conducted experiments using two publicly accessible datasets to validate their approach. Wise et al. [11] developed an attention-enhanced convolutional neural network specifically designed for accurately classifying a wide range of 19 orchestral instruments, and they substantiated the effectiveness of this technique through experimental evaluations.

## III. INSTRUMENT FEATURE EXTRACTION

Different musical instruments emit distinct vibration frequencies, which, in turn, yield different sounds. Timbre represents the subjective, natural perception of these vibrations by humans. To achieve the automatic classification of various musical instruments, it is essential to extract features that reflect the timbral characteristics of the musical instrument sound signal. Currently, in the realm of speech signal recognition, a predominant reliance on time-domain and frequency-domain features is observed. These feature types are relatively straightforward and easy to analysis. However, they fall short in adequately characterizing timbre, resulting in suboptimal classification performance. This paper primarily focuses on the analysis of two categories of features derived from the human auditory system for musical instrument feature extraction.

MFCC is a feature obtained by simulating the auditory characteristics of the human ear [12], which has a good performance in speech recognition [13]. It is extracted in the following way.

*1)* Pre-emphasis, frame-by-frame, and windowing operations are performed on the original instrument audio signal to get pre-processed signal $x_i(m)$.

*2)* A fast Fourier transform (FFT) is performed on $x_i(m)$, and

$$X(i,t) = FFT[x_i(m)] \tag{1}$$

is obtained.

*3)* The energy spectrum is calculated:
$$E(i,t) = [x_i(m)]^2. \tag{2}$$

*4)* A filtering operation is performed on $E(i,t)$. The spectrum energy is obtained through a group of triangular filters:
$$S(i,m) = \sum_{t-0}^{N-1} E(i,t)H_m(t), \tag{3}$$

where, $H_m(t)$ represents the transfer function of the triangular filter.

*5)* The final MFCC is obtained by discrete cosine transform (DCT):

$$X(i) = \sqrt{\frac{2}{M}} \sum_{m=0}^{M-1} lg[S(i,m)] \cos\left[\frac{\pi(2m-1)}{2M}\right], \tag{4}$$

where, $M$ is the filter order.

PLPC is also a feature that mimics the auditory characteristics of the human ear [14] and is extracted as follows:

*1)* FFT is also performed on $x_i(m)$ to convert it to the frequency domain to obtain $X(i,t)$. $E(i,t)$ is calculated.

*2)* $E(i,t)$ is converted to the Bark domain:

$$\varphi(f) = 6\ln\left\{\frac{f}{600} + \left[\left(\frac{f}{600}\right)^2 + 1\right]^{0.5}\right\}, \tag{5}$$

where, $\varphi$ is the value of the signal frequency in the Bark domain and $f$ is the angular frequency.

*3)* The critical bands in the Bark domain are considered as a group of filters, and the function of each filter is:

$$\varphi(Z - Z_0(k)) = \begin{cases} 0, Z - Z_0(k) < -1.3 \\ 10^{Z-Z_0(k)+0.5}, -1.3 \leq Z - Z_0(k) \leq -0.5 \\ 1, -0.5 < Z - Z_0(k) < 0.5 \\ 10^{-2.5(Z-Z_0(k)-0.5)}, 0.5 \leq Z - Z_0(k) \leq 2.5 \\ 0, 2.5 < Z - Z_0(k) \end{cases} \tag{6}$$

The output of each critical band is obtained:

$$\theta(k) = \sum_{Z-Z_0(k)=-1.3}^{2.5} E(\varphi - \varphi_i)\,\varphi(Z - Z_0(k)). \tag{7}$$

*4)* Equal loudness curve pre-emphasis is carried out on $\theta(k)$:

$$\Gamma(k) = D[f_0(k)]\theta(k), \tag{8}$$

where, $f_0(k)$ corresponds to central frequency $Z_0(k)$ of each filter, $Z_0(k)$ refers to the equal loudness curve,

$$D[f_0(k)] = \frac{16\pi^4 f_0^4(k)(4\pi^2 f_0^2(k) + 5.68\times10^7)}{(4\pi^2 f_0^2(k) + 6.3\times10^6)^2 \times (4\pi^2 f_0^2(k) + 3.8\times10^8)}. \tag{9}$$

*5)* Finally, the simulation of human ear hearing is realized by converting from intensity to loudness:
$$\Psi(k) = \Gamma(k)^{0.33}. \tag{10}$$

In the context of speech recognition tasks, MFCC is usually set as 13 dimensions along with their corresponding first-order and second-order difference coefficients, resulting in a total of 39 dimensions. Additionally, PLPC usually employs 24 filters to yield 24 dimensions. If both the 39-dimensional MFCC and the 24-dimensional PLPC features are utilized for automatic multi-instrument classification, the cumulative dimensionality reaches 63 dimensions. To mitigate the potential complexity and computational load associated with this higher dimensionality, this paper uses an entropy weight method [15] to optimize these features.

Information entropy is related to the probability of an event occurring. When it is applied to feature optimization, the greater the information entropy, the more disordered the distribution of that feature value is. For a $T_n \times T_m$ original feature set, it is assumed that the feature vector corresponding to the $t_m$-th feature item is $Z_{t_m}$, the process of feature optimization is presented in Fig. 1.

Fig. 1.   The feature optimization process based on entropy weight method.

The specific calculation process is as follows:

*1)* Standardized eigenvalue:

$$Z'_{t_n,t_m} = \frac{z_{t_n,t_m} - min(Z_{t_m})}{max(Z_{t_m}) - min(Z_{t_m})}. \tag{11}$$

*2)*  The information entropy is calculated:

$$E(t_m) = -\sum_{t_n=1}^{T_n} p_x(t_n, t_m) \ln p_x(t_n, t_m), \tag{12}$$

where, $p_x(t_n, t_m)$ refers to the percentage of different eigenvalues in the corresponding term, $p_x(t_n, t_m) = \frac{X'(t_n, t_m)}{P_z(t_m)}$, $P_z(t_m) = \sum_{t_n=1}^{T_n} X'(t_n, t_m)$.

*3)*  The feature term weight is calculated:

$$W_{t_m} = \frac{1 - E(t_m)}{T_m - \sum_{t_m=1}^{T_m} E(t_m)}. \tag{13}$$

*4)* A threshold is set, and feature items with weights greater than the threshold are retained as subsequent features used for automatic multi-instrument classification.

## IV.   AUTOMATIC CLASSIFICATION METHOD BASED ON NEURAL NETWORKS

Back-propagation neural network (BPNN) has excellent performance and is most widely used in automatic classification tasks [16]. Therefore, in this paper, BPNN is chosen as a classifier for automatic classification of multi-instrument. A simple BPNN structure is presented in Fig. 2.



Fig. 2.   The BPNN structure.

It is hypothesized that the number of nodes in the input layer, hidden layer, and output layer of BPNN is $m - p - q$. The output of the hidden layer during training can be written as:

$$h_j = f\left(\sum_{i=1}^{m} w_{ij}^h i_i - b_j\right). \tag{14}$$

The output of the output layer can be written as:

$$o_k = f\left(\sum_{j=1}^{p} w_{jk}^o h_j - b_k\right). \tag{15}$$

where, $w$ and $b$ denote the weight and threshold of each layer.

The global error function can be written as:

$$E = \frac{1}{2}\sum_{k=1}^{q}(y_k - o_k)^2, \tag{16}$$

where, $y_k$ is the desired output. The BPNN continuously corrects the weight threshold according to the error until the error meets the requirements. However, the performance of BPNN is easily affected by the initial weight threshold, so it can be tempting to become trapped in the local optimum. In order to solve this problem, this paper uses the sparrow search algorithm (SSA) [17] to optimize the BPNN.

SSA is a sparrow-inspired algorithm that utilizes the foraging behavior of sparrows, which exhibits excellent global optimization performance and is used to calculate the optimal initial weight thresholds for BPNN. During the foraging process, some sparrows are responsible for finding the foraging area and direction, the rest of the sparrows feed, and some sparrows will sound an alarm when danger is detected. The population is categorized into discoverers, followers, and alerts, and their positions are updated as follows.

*1)* Discoverer:

$$X_{i,j}^{t+1} = \begin{cases} X_{i,j}^t \cdot exp\left(-\frac{i}{\alpha \cdot t_{max}}\right), R < A_{ST} \\ X_{i,j}^t + QL, R \geq A_{ST} \end{cases}, \tag{17}$$

where, $\alpha$ is a random number in $[0,1]$, $R$ is an alert value in $(0,1]$, and $A_{ST}$ is the threshold at which the shift occurs, $Q$ is a random number obeying a normal distribution, $L$ is a $1 \times d$ matrix, $R < A_{ST}$ means that the area is dangerous and the finder needs to move to a safe area, and $R \geq A_{ST}$ means that the area is safe and sparrows can look for food.

*3)* Follower:

$$X_{i,j}^{t+1} = \begin{cases} Q \cdot exp\left(\frac{X_{worst}^t - X_{i,j}^t}{i^2}\right), i > \frac{n}{2} \\ X_b^{t+1} + \left|X_{i,j}^t - X_b^{t+1}\right| \cdot A^+ \cdot L, i \le \frac{n}{2} \end{cases}, \quad (18)$$

where, $X_{worst}$ is the worst position of the individual, $X_b$ is the best position of the discoverer, $A^+ = A^T(AA^T)^{-1}$, $A$ is a $1 \times d$ matrix, randomly assigned as -1 or 1.

*4)* Watcher:

$$X_{i,j}^{t+1} = \begin{cases} X_{best}^t + \beta \cdot \left|X_{i,j}^t - X_{best}^t\right|, f_i > f_g \\ X_{i,j}^t + K \cdot \left(\frac{\left|X_{worst}^t - X_{i,j}^t\right|}{f_i - f_w + \varepsilon}\right), f_i = f_g \end{cases}, \quad (19)$$

where, $X_{best}$ is the global optimal position, $\beta$ is the step control parameter, $f_i$ refers to the fitness of the i-th sparrow, $f_g$ and $f_w$ are the current best and worst fitness, $K$ is a random number in (-1,1), and $\varepsilon$ is a constant.

SSA finds the optimal weight threshold of BPNN by constantly updating the three positions. To improve the population diversity, this paper uses cubic mapping [18] to obtain the initialized population:

$$x_{n+1} = \rho x_n(1 - x_n^2), \quad (20)$$

where, $\rho$ is the control parameter. The SSA-BPNN algorithm is used to realize the classification of different musical instruments, and the flow chart is presented in Fig. 3.



Fig. 3. The SSA-BPNN algorithm-based automatic multi-instrument classification method.

## V. RESULTS AND ANALYSIS

### A. Experimental Setup

The algorithm was developed and trained in the MATLAB 2018b environment. The dataset used for the experiment was the IRMAS dataset [19], comprising a total of 6,705 WAV audio files. All audio files were in 16-bit stereo format and had a sampling rate of 44.1 kHz. A single performance clip of ten instruments, including cello, clarinet, and others, each with a duration of 3 s, was selected. The training set and test set configurations are detailed in Table I.

TABLE I. EXPERIMENTAL DATA SETS

| Musical Instrument | Number in the Training Set | Number in the Test Set |
|---|---|---|
| Flute | 451 | 163 |
| Organ | 682 | 361 |
| Piano | 721 | 995 |
| Trumpet | 577 | 167 |
| Cello | 388 | 111 |
| Clarinet | 505 | 62 |
| Electric guitar | 760 | 942 |
| Violin | 580 | 211 |
| Saxophone | 626 | 326 |
| Acoustic guitar | 637 | 535 |

The classification results were assessed using the precision (P), recall rate (R), and F1 value, calculated by:

$$P = \frac{TP}{TP+FP}, \quad (21)$$

$$R = \frac{TP}{TP+FN}, \quad (22)$$

$$F_1 = \frac{2PR}{P+R}, \quad (23)$$

where:

$TP$: number of samples retrieved and belonging to the positive category,

$FP$: number of samples retrieved but belonging to the negative category,

$FN$: number of samples not retrieved but belonging to the positive category.

In automatic multi-instrument classification, the macroscopic values of P and R need to be considered:

$$P_{marco} = \frac{1}{|L|}\sum_{l=1}^{L} P_l, \quad (24)$$

$$R_{marco} = \frac{1}{|L|}\sum_{l=1}^{L} R_l, \quad (25)$$

where, $L$ refers to the number of labels, $P_l$ and $R_l$ are the corresponding P and R values calculated for each label $l$. By further calculation, the corresponding macro F1 value can be obtained.

### B. Result Analysis

Firstly, the effect of feature selection on the multi-instrument classification results was analyzed. In the case of using the SSA-BPNN model as a classifier but changing the algorithm's feature input, the classification results were compared, as shown in Table II.

TABLE II.     EFFECT OF FEATURE SELECTION ON MULTI-INSTRUMENT CLASSIFICATION RESULTS

| Musical Instrument | | Flute | Organ | Piano | Trumpet | Cello | Clarinet | Electric Guitar | Violin | Saxophone | Acoustic Guitar | Macro-value |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MFCC | P | 0.69 | 0.61 | 0.55 | 0.66 | 0.69 | 0.61 | 0.62 | 0.63 | 0.59 | 0.71 | 0.64 |
| | R | 0.66 | 0.64 | 0.65 | 0.62 | 0.66 | 0.71 | 0.65 | 0.61 | 0.71 | 0.61 | 0.65 |
| | F1 | 0.67 | 0.62 | 0.60 | 0.64 | 0.67 | 0.66 | 0.63 | 0.62 | 0.64 | 0.66 | 0.64 |
| PLPC | P | 0.71 | 0.65 | 0.66 | 0.65 | 0.68 | 0.65 | 0.62 | 0.65 | 0.67 | 0.68 | 0.66 |
| | R | 0.65 | 0.67 | 0.66 | 0.64 | 0.66 | 0.68 | 0.71 | 0.68 | 0.72 | 0.66 | 0.67 |
| | F1 | 0.68 | 0.66 | 0.66 | 0.64 | 0.67 | 0.66 | 0.66 | 0.66 | 0.69 | 0.67 | 0.67 |
| MFCC+PLPC | P | 0.71 | 0.65 | 0.65 | 0.68 | 0.74 | 0.68 | 0.71 | 0.61 | 0.68 | 0.68 | 0.68 |
| | R | 0.65 | 0.61 | 0.65 | 0.65 | 0.71 | 0.71 | 0.65 | 0.61 | 0.75 | 0.74 | 0.67 |
| | F1 | 0.68 | 0.63 | 0.65 | 0.66 | 0.72 | 0.69 | 0.68 | 0.61 | 0.71 | 0.71 | 0.68 |
| MFCC + PLPC (entropy weight method) | P | 0.75 | 0.66 | 0.68 | 0.79 | 0.81 | 0.72 | 0.71 | 0.64 | 0.71 | 0.71 | 0.72 |
| | R | 0.68 | 0.62 | 0.68 | 0.71 | 0.75 | 0.75 | 0.75 | 0.66 | 0.78 | 0.75 | 0.71 |
| | F1 | 0.71 | 0.64 | 0.68 | 0.75 | 0.78 | 0.73 | 0.73 | 0.65 | 0.74 | 0.73 | 0.71 |

According to the data in Table II, when MFCC was used as the feature, the multi-instrument classification yielded a P value of 0.71, an R value of 0.65, and an F1 value of 0.64. In contrast, when PLPC was employed as the feature, the multi-instrument classification produced a P value of 0.66, which was 0.02 larger than that when using MFCC. The R value was 0.67, which was 0.02 larger than that when using MFCC. The F1 value was 0.67, which was 0.03 larger than that when using MFCC. These results suggested that PLPC was more effective in the context of automatic multi-instrument classification when compared to MFCC. It can be attributed to the fact that PLPC is better at simulating the human ear's perception of sound, making it closer to actual music perception.

Subsequently, both MFCC and PLPC were simultaneously input into the SSA-BPNN model, resulting in a P value of 0.68, which showed an increase of 0.02 compared to inputting PLPC alone. The R value was 0.67, the same as when inputting PLPC alone, and the F1 value reached 0.68, which was an increase of 0.01 compared to inputting PLPC alone. These outcomes suggested that the enhancement in classification performance achieved by inputting both sets of features into the SSA-BPNN model simultaneously was not substantial. It may be because the large number of features affects the algorithm classification performance.

In the final step, the MFCC+PLPC features, selected through the entropy weight method, were employed. This configuration yielded a P value of 0.72, an R value of 0.71, and an F1 value of 0.71. These results surpassed the previous three feature input combinations, underscoring the effectiveness of utilizing the entropy weighting method for feature optimization. This approach enhances the automatic classification capability of the SSA-BPNN method for multi-instrumental instruments while simultaneously reducing the number of feature dimensions.

Then, also using MFCC+PLPC after the entropy weighting method selection as features, the effect of different classifiers on the multi-instrument classification results was compared. The SSA-BPNN algorithm was compared with other classifiers, including:

*1)* Support vector machine (SVM) [20],

*2)* XGBoost [21],

*3)* Traditional BPNN [22],

*4)* BPNN optimized using genetic algorithms (GA): GA-BPNN [23],

*5)* BPNN optimized using particle swarm algorithm (PSO): PSO-BPNN [24].

The results are presented in Fig. 4.



Fig. 4.    The effect of the classifier on the automatic multi-instrument classification results.

From Fig. 4, it is evident that the BPNN algorithm exhibited superior classification performance compared to the SVM and XGBoost algorithms. The BPNN algorithm achieved a P value of 0.65, which were improved by 0.12 compared to the SVM algorithm and 0.04 compared to the XGBoost algorithm. The R value for the BPNN algorithm was 0.61, representing a 0.09 increase compared to the SVM algorithm and a 0.02 increase compared to the XGBoost algorithm. The F1 value for the BPNN algorithm was 0.63, indicating an improvement of 0.11 compared to the SVM algorithm and an increase of 0.11 compared to the XGBoost algorithm. These

results underscored the advantages of using BPNN as a classifier. In the comparison of different weight threshold optimization methods, the GA-BPNN algorithm showed P value, R value, and F1 value scores below 0.7 for automatic multi-instrument classification, which did not demonstrate substantial improvement over the BPNN algorithm. The PSO-BPNN algorithm achieved only a P value of 0.71, an R value of 0.68, and an F1 value of 0.69. Compared to the PSO-BPNN algorithm, the SSA-BPNN algorithm showed an increase of 0.01 in P value, 0.03 in R value, and 0.02 in F1 value. This result highlighted the reliability of the classifier developed in this paper for automatic multi-instrument classification.

## VI. Discussion

The automatic classification of multiple musical instruments is an important area of research in audio signal processing [25], and it also has significant implications for both theoretical and practical applications. Through the automatic classification of multiple musical instruments, it can provide strong support for MIR [26], helping users conveniently find music they are interested in. In music composition, it can be used to automatically separate different instrument tracks, providing flexible post-production techniques for music producers. In cultural preservation, it enables the analysis of multiple instruments in traditional music [27], contributing to better protection and inheritance of musical culture. However, music performance often involves a wide variety of complex instruments, which poses significant challenges for automatic classification. Further research is needed to enhance the effectiveness of multi-instrument automatic classification.

In current research on the automatic classification of multiple musical instruments, improving classification accuracy mainly relies on optimizing feature extraction and classification algorithms. In terms of feature extraction, this paper utilized entropy weighting to optimize MFCC and PLPC features, reducing dimensionality while enhancing accuracy. Regarding the classification algorithm, SSA was employed to optimize BPNN parameters for better performance. The experiments conducted on the IRMAS dataset revealed that the feature selection optimization method, based on entropy weighting, effectively enhanced the accuracy of multi-instrument automatic classification. The obtained P value was 0.72, the R value was 0.71, and the F1 value was 0.71, which were higher than those achieved by other feature combinations. This result indicated that the features selected through entropy weighting could better represent the characteristics of different instruments, thereby improving the discrimination effect of the SSA-BPNN method for various instruments. The analysis of the classifier revealed that the SSA demonstrated excellent performance in optimizing the parameters of BPNN and it could enhance classification accuracy and obtain results better than the other classifiers. This makes it applicable in practical scenarios.

## VII. Conclusion

This paper proposed an approach for the automatic classification of multiple instruments. The MFCC and PLPC features were optimized using the entropy weighting method. An SSA-BPNN method was designed as the classifier. Through experiments, it was found that the features optimized by the entropy weighting method delivered optimal performance in the automatic classification of multiple instruments, with a P value of 0.72, an R value of 0.71, and an F1 value of 0.71, outperforming alternative methods like the SVM algorithm.

This method can accurately classify multiple musical instruments and can be further applied in practical music data processing.

The findings of this study demonstrate that the SSA-BPNN method, as designed, can achieve a relatively high level of accuracy in automatically classifying multiple instruments. This provides valuable theoretical insights for optimizing features, refining neural network algorithms, and improving parameter optimization methods. Moreover, it opens up new avenues for future research on automatic classification of multiple instruments. From a practical standpoint, the proposed method holds potential for application in real-world scenarios involving music data processing and even speech data processing.

However, this study also has certain limitations. For instance, the effectiveness of the proposed method has not been validated on a larger dataset, and there is potential for further optimization in feature selection. In future research, it is imperative to conduct additional investigations into feature combinations and optimization methods, explore strategies to enhance the classification accuracy of the model, and examine the applicability of this approach in domains such as artist classification and music genre categorization.

## References

[1] S. Rajesh, and N. J. Nalini, "Musical instrument emotion recognition using deep recurrent neural network," Proc. Comput. Sci., vol. 167, pp. 16-25, January 2020.

[2] C. E. Cella, "Music information retrieval and contemporary classical music: a successful failure," Trans. Int. Soc. Music Inform. Retriev., vol. 3, pp. 126-136, September 2020.

[3] X. Fu, H. Deng, and J. Hu, "Automatic label calibration for singing annotation using fully convolutional neural network," IEEJ T. Electr. Electr., vol. 18, pp. 945-952, April 2023.

[4] D. Martin-Gutierrez, G. H. Penaloza, A. Belmonte-Hernandez, and F. Alvarez Garcia, "A multimodal end-to-end deep learning architecture for music popularity prediction," IEEE Access, vol. 8, pp. 39361-39374, February 2020.

[5] M. Voetter, M. Mayerl, G. Specht, and E. Zangerle, "HSP datasets: insights on song popularity prediction," Int. J. Semant. Comput., vol. 16, pp. 233-255, June 2022.

[6] S. Rajesh, and N. J. Nalini, "Combined evidence of MFCC and CRP features using machine learning algorithms for singer identification," Int. J. Pattern Recogn., vol. 35, pp. 2158001.1-2158001.21, July 2020.

[7] S. Kooshan, H. Fard, and R. M. Toroghi, "Singer identification by vocal parts detection and singer classification using LSTM neural networks," 2019 4th Int. Conf. on Pattern Recognition and Image Analysis (IPRIA), pp. 246-250, March 2019.

[8] W. Feng, J. Liu, T. Li, Z. Yang, and D. Wu, "FAC: a music recommendation model based on fusing audio and chord features (115)," Int. J. Softw. Eng. Know., vol. 32, pp. 1753-1770, October 2022.

[9] A. Elbir and N. Aydin, "Music genre classification and music recommendation by using deep learning," Electron. Lett., vol. 56, pp. 627-629, June 2020.

[10] D. Chatterjee, A. Dutta, D. Sil, and A. Chandra, "Deep single shot musical instrument identification using scalograms," 2023 Int. Conf. on

Artificial Intelligence in Information and Communication (ICAIIC), pp. 386-389, August 2023.

[11] A. Wise, A. S. Maida, and A. Kumar, "Attention augmented CNNs for musical instrument identification," 2021 29th Eur. Signal Process. Conf. (EUSIPCO), pp. 376-380, August 2021.

[12] S. K. Mahanta, A. F. U. R. Khilji, and P. Pakray, "Deep neural network for musical instrument recognition using MFCCs," Comput. Sist., vol. 25, pp. 351-360, May 2021.

[13] E. Nakamura, Y. Kageyama, and S. Hirose, "LSTM-based japanese speaker identification using an omnidirectional camera and voice information," IEEJ T. Electr. Electr., vol. 17, pp. 674-684, January 2022.

[14] V. V. Yerigeri, and L. K. Ragha, "Speech stress recognition using semi-eager learning," Cogn. Syst. Res., vol. 65, pp. 79-97, January 2021.

[15] T. Xu, C. Qin, H. Zhang, Y. Qu, and W. Fang, "Study on petroleum standard attention index calculation based on the entropy weight method," IOP Conf. Ser.: Earth Environ. Sci., vol. 514, pp. 1-7, July 2020.

[16] Y. Liu, L. Nie, R. Dong, and G. Chen, "BP neural network-Kalman filter fusion method for unmanned aerial vehicle target tracking," Proc. Inst. Mech. Eng. C: J. Mec. Eng. Sci., vol. 237, pp. 4203-4212, July 2023.

[17] G. F. Fan, Y. Li, X. Y. Zhang, Y. H. Yeh, and W. C. Hong, "Short-term load forecasting based on a generalized regression neural network optimized by an improved sparrow search algorithm using the empirical wavelet decomposition method," Energy Sci. Eng., vol. 11, pp. 2444-2468, April 2023.

[18] A. Zahid, and M. Arshad, "An innovative design of substitution-boxes using cubic polynomial mapping," Symmetry, vol. 11, pp. 1-10, March 2019.

[19] J. J. Bosch, J. Janer, F. Fuhrmann, and P. Herrera, "A comparison of sound segregation techniques for predominant instrument recognition in musical audio signals," Proc. ISMIR, pp. 559-564, January 2012.

[20] F. Camastra, V. Capone, A. Ciaramella, A. Riccio, and A. Staiano, "Prediction of environmental missing data time series by support vector machine regression and correlation dimension estimation," Environ. Modell. Softw., vol. 150, pp. 1-7, April 2022.

[21] Y. Wu, Q. Zhang, Y. Hu, S. W. Ko, X. Zhang, H. Zhu, J. Liu, and S. Li, "Novel binary logistic regression model based on feature transformation of XGBoost for type 2 Diabetes Mellitus prediction in healthcare systems," Future Gener. Comp. Sy., vol. 129, pp. 1-12, November 2022.

[22] M. Wei, X. Hu, and H. Yuan, "Residual displacement estimation of the bilinear SDOF systems under the near-fault ground motions using the BP neural network," Adv. Struct. Eng., vol. 25, pp. 552-571, December 2021.

[23] J. He, X. Li, and Y. Zhao, "The fault diagnosis of diesel fuel supply system based on BP neural network optimized by genetic algorithm," J. Phys.: Conf. Ser., vol. 1732, pp. 1-6, January 2021.

[24] Y. Wang, and Y. Zhao, "Predicting bedrock depth under asphalt pavement through a data-driven method based on particle swarm optimization-back propagation neural network," Constr. Build. Mater., vol. 354, pp. 1-16, November 2022.

[25] S. R. Gulhane, S. D. Shirbahadurkar, and S. S. Badhe, "Indian classical musical instrument classification using timbral features," Concurr. Comp. Pract. E., vol. 33, pp. e6414.1-e6418.19, May 2021.

[26] Y. H. Chen, M. Ceccarelli, and H. S. Yan, "A historical study and mechanical classification of ancient music-playing automata," Mech. Mach. Theory, vol. 121, pp. 273-285, March 2018.

[27] S. R. Chaudhary and S. N. Kakarwal, "Various approaches in musical instrument identification: a review," Int. J. Appl. Evol. Comput., vol. 10, pp. 1-7, April 2019.

# HarborSync: An Advanced Energy-efficient Clustering-based Algorithm for Wireless Sensor Networks to Optimize Aggregation and Congestion Control

Ibrahim Aqeel

College of Computer Science & IT, Jazan University, Jazan, Saudi Arabia

*Abstract*—In the ever-evolving landscape of Wireless Sensor Networks (WSNs), the demand for cutting-edge algorithms has never been more critical. This paper proposes an algorithm, HarborSync, to improve stability, energy efficiency, durability, and congestion control in WSN. While selecting cluster heads and backup nodes, HarborSync applies the Optimised Stable Clustering Algorithm (OSCA) and the Weighted Clustering Algorithm (WCA). This fresh method puts the groundwork for better performance by acquainting techniques to intentionally postpone changes in cluster heads and computing priorities. Using the innovative Cluster-based Aggregation and Congestion Control (CACC) features, HarborSync provides enhanced routing, adaptive reconfiguration, efficient aggregation techniques, and dynamic congestion monitoring. Among HarborSync's strengths, stability bears out with a 90% rating, surpassing those of LEACH (78%), LEACH-C (82%), TEEN (88%), and PEGASIS (76%). When it comes to durability, HarborSync scores 88% better than LEACH (75%), LEACH-C (80%), TEEN (85%), and PEGASIS (72%). The HarborSync score is 3.85% for congestion control compared to LEACH and LEACH-C, managing 5.22%, TEEN accomplishing 4.98%, and PEGASIS with 7.32%. Regarding adaptability, HarborSync showcases its versatility, earning an 85% rating, surpassing LEACH (72%), competes with LEACH-C (78%), equals TEEN (90%), and outperforms PEGASIS (68%). In the critical realm of packet loss management, HarborSync demonstrates efficiency with a reduced rate of 6.179%. Therefore, it outperforms LEACH (7.811%), LEACH-C (6.897%), TEEN (4.953%), and PEGASIS (7.973%).

*Keywords*—*Clustering; congestion control; cluster head selection; energy-efficient clustering; wireless sensor networks; energy optimization*

## I. INTRODUCTION

To monitor and sense faraway places, WSNs collect data from many tiny, inexpensive sensors [1] and send it to a central station. WSNs are versatile and affordable, making them useful in many fields, including healthcare, emergency response [2], weather prediction, and surveillance missions. On top of that, these networks have the capability to build ad hoc networks, which allow them to operate autonomously in challenging or dangerous environments when human intervention is not an option. The operating longevity of WSNs depends on battery life [3], which can be difficult, if not impossible, to replenish, making energy efficiency a significant concern. Consequently,

optimizing network functionality necessitates the creation of algorithms that consume less energy. Various algorithms that minimize energy consumption have been suggested for WSNs in the past few years [3].

When neighboring nodes detect the same or comparable events, clustering becomes an effective method for arranging ad hoc sensors. Congestion and data collisions result from the network's energy being quickly depleted due to individual communication between each node and the base station [4]. To solve this problem, clustering organizes sensor nodes into smaller groups, each with its designated cluster head (CH). Each node uses short-distance transmission to send sensed data to its corresponding CH, and then each CH uses long-distance transmission to aggregate and send the aggregated data to the base station. As a result, CHs use more power while sending messages than other cluster members [5]. It has been suggested that clusters hold CH elections periodically to reduce energy imbalances. To improve the network's lifetime, it is essential to determine the ideal number and size of clusters. As data is transferred from cluster members to cluster leaders [6], an excessive amount of energy is consumed when the number of clusters is minimal. On the other side, many nodes have to use long-distance transmission to talk to the base station when there are many clusters since so many cluster leaders were elected. To maximize energy usage across the network, it is necessary to strike a balance between these parameters [7]. For optimal intra-cluster performance, it is critical to distribute cluster heads evenly over the network. When cluster heads are chosen too closely, clusters don't develop evenly, making efficient clustering difficult since cluster sizes are too small or too big. Overhearing signals and wasting energy can happen in any case. Modern technological landscapes bet on WSNs for environmental monitoring, healthcare, industrial automation, and smart cities, among many others [8]. As these networks become increasingly constitutional to our unified society, there is a compacting need for sophisticated algorithms to heighten their functionality [9]. In response to the increasing challenges run into by WSNs, this research infixes "HarborSync," a novel algorithm. In the era of WSNs, cutting-edge solutions are crucial since these networks must voyage the challenges of guaranteeing constant data transfer [10], minimizing energy consumption, and maintaining network resilience over time.

By strategically expanding upon the Optimised Stable Clustering Algorithm (OSCA), HarborSync uses the Weighted Clustering Algorithm (WCA) to choose accompaniment nodes and cluster heads carefully. This one-of-a-kind combination premises features including priority computing, purposeful delay in cluster head transitions, and the incorporation of Cluster-based Aggregation and Congestion Control (CACC) [11], laying the groundwork for a path-breaking approach. At last, we have an all-inclusive solution that surmounts all subsisting algorithms on various metrics, including stability, energy efficiency, durability, adaptability, congestion control, and packet loss management. Section I of the paper introduces the study subject and provides background information on WSN and clustering [12]. A thorough understanding of the HarborSync algorithm is the goal of the paper's organization. The second section lays the groundwork for the proposed technique and examines the current body of knowledge through a thorough literature study. The new method, HarborSync, is described in full in Section III, along with its components and the design rationale of the proposed algorithms. To demonstrate how HarborSync outperforms well-known algorithms like LEACH, LEACH-C, TEEN, and PEGASIS, Section IV presents the results showing how well the algorithm performs when tested with different evaluating settings. Section V presents the discussion. Key findings are summarised in the conclusion, and additional study and development possibilities are suggested in Sections VI and VII in the future scope section.

The following are the contributing points that have been addressed in the study based on the proposed algorithm:

*1)* In this study, a new technique called HarborSync is introduced and proposed; it is particularly made for Wireless Sensor Networks (WSNs). The HarborSync uses the Weighted Clustering Algorithm (WCA) features and the Optimised Stable Clustering Algorithm (OSCA) to choose backup nodes and cluster chiefs.

*2)* We suggested HarborSync offers Cluster-based Aggregation and Congestion Control (CACC) capabilities to enhance the effectiveness and dependability of data transmission in the network, demonstrating creativity in resolving data aggregation and congestion problems in WSNs.

*3)* A thorough performance comparison between HarborSync and well-known clustering techniques like LEACH, LEACH-C, TEEN, and PEGASIS is included in the study. This comparative study emphasizes the algorithm's virtues, which shows how it outperforms current techniques in terms of packet loss management, stability, durability, flexibility, and congestion control.

*4)* One of HarborSync's main benefits is its capacitance to lower power consumption, which is essential for wireless sensor networks. The article hashes out the algorithm's manifested capacity to save power, which bestows the reliability and endurance of WSNs.

*5)* Panoptic testing equates HarborSync to popular clustering algorithms on many parameters. It raises HarborSync's functionality testing in many scenarios.

## II. RELATED LITERATURE

One of the first cluster-based routing concepts for WSNs was LEACH, which was proposed by Heinzelman et al. [3]. LEACH uses a random rotation of the cluster heads to ensure that all sensor nodes use the same amount of power. Because it is not controlled by a single entity, its decentralized nature makes it ideal for networks with many nodes. But, because LEACH is inherently random, specific nodes may experience early energy depletion or an unbalanced energy distribution. In 2003, the same authors introduced LEACH-C [4] as an improvement to LEACH. Resolving some of LEACH's shortcomings, it presents a centralized mechanism for selecting cluster heads. Using parameters like residual energy and distance from the base station, LEACH-C uses a base station to identify cluster heads. This centralized strategy aims to increase the network's lifetime and decrease its energy consumption. Still, problems with centralization and base station connectivity pose a continuing threat. This study offer a new threshold function for cluster head selection that optimizes the LEACH protocol, resulting in an energy efficient clustering method for FANETS [13]. The results from the MATLAB experiments show that the new protocol is more energy-efficient than the current LEACH and Centralized Low-Energy Adaptive Clustering Hierarchy Protocol. It also has a higher packet delivery rate and a lower First Node Death (FND). A referenced study [5] assessed the present clustering routing protocols, categorizing these algorithms into 2 primary types of roting techniques namely data transmission and cluster-construction. The review considered sixteen well-established clustering methods, excluding newer approaches like fuzzy and evolutionary-based methods [14].

In the context of Wireless Sensor Networks (WSNs), this study probed node clustering methods founded on fuzzy modeling. The basal concentrate was on assessing their benefits and drawbacks. Classifying clustering algorithms as fuzzy or hybrid fuzzy-based was one panorama of the inquiry. Diverse methodologies were engaged in a different study [6] to investigate cluster-based routing schemes. Canvassing these techniques fractioned into block, chain, and grid-based methods showcased their benefits and drawbacks. Cluster stability, scalability, energy economy, and delivery time were the valuing retainers [15].

Ramping upon this basis, more research [7] categorized cluster-based routing algorithms as block, grid, or chain-based by probing clustering protocols. The comparative valuations of stream feelers considered factors such as algorithmic complexity, load balancing, cluster stability, delivery time, energy cognizance, and load sensitivity [16]. A paper [9] dealt with the problem of classifying several WSN clustering techniques into heterogeneous and homogeneous networks. This study essayed the advantages and disadvantages of each protocol while describing the network node and resource capacity. The equivalence research admitted Cluster Heads (CHs), complexity, number of clusters, clustering items, and inter-cluster communication [3].

Furthermore, unequal clustering techniques were examined in [10] based on their attributes and objectives. The comparability focused on the clustering process and clustering

attributes, forking the techniques into three categories: deterministic, preset, and probabilistic methods. Legion methods were also simulated in order to gauge their energy usage and service life. Heterogeneous and homogeneous networks (Energy-Efficient Stable LEACH) [11] are variations of the LEACH protocol contrived to increase the energy efficiency of the network. The particle swarm optimization (PSO) technique was used in this variation. In the ESO-LEACH model, each node utilizes a probability descent from the ESO algorithm and considers its energy level when selecting a Cluster Head (CH). Based on the current energy levels and node distances, this method depends on a CH with enough energy to subsist the whole round. One substantial drawback of ESO-LEACH is its computational cost. It surpasses that of the original LEACH protocol. It becomes problematic to implement on devices with limited resources, and the protocol would have trouble adjusting to changing network conditions, which would require recalculating the clustering structure from the ground up.

A load-balancing technique was developed in separate research [12] to improve the efficiency of 5G Local Home Networks (5GLHNs). The CFPSO (Cell Attachment using Particle Swarm Optimization) technique was utilized in this process for cell attachment. A separate inquiry examined techniques and approaches for achieving precise time synchronization in femtocell networks [17]. It proposed an intra-cluster synchronization mechanism to improve the accuracy of synchronization. The proposed technique was subjected to empirical testing to assess its consumption of resources and its security features. In addition, another research group has devised an energy-efficient approach for selecting CH (Cluster Head) that considers many criteria, including remaining energy, distance, and node density, utilizing Particle Swarm Optimization (PSO) [18]. Nevertheless, this approach fails to consider the clustering procedure, leading to significant energy inefficiency throughout the network, and fails to acknowledge the creation of clusters. Table I reviews the state-of-the-art literature, showing its advantages and disadvantages.

TABLE I. RELATED LITERATURE

| Ref | Advantages | Gaps |
|---|---|---|
| LEACH [3] | -The decentralized nature is suitable for networks with many nodes. <br> - Random rotation of cluster heads for energy balance. | - Inherently random, leading to early energy depletion or unbalanced energy distribution. |
| LEACH-C [4] | - Centralized mechanism improves energy efficiency. <br> - Selection of cluster heads based on residual energy and distance from the base station. | - Centralization and base station connectivity issues pose threats. |
| EE-LEACH [5] | - Categorizes clustering algorithms into data-transmission and cluster-construction routing techniques. <br> - -Energy efficient clustering for FANETS | - Excludes newer approaches like fuzzy and evolutionary-based methods. |
| Fuzy based clustering [6] | - Focuses on merits and limitations of fuzzy modelling-based node clustering methods. <br> - Classifies fuzzy and hybrid fuzzy-based clustering methods. | - Limited to fuzzy modeling approaches, excluding other clustering techniques. |
| heterogeneous and homogeneous clustering [9] | - Classifies cluster-based routing techniques into block, grid, and chain-based methods. <br> - Comparative evaluations cover delivery delay, energy consumption, load balance, cluster strength, and complexity of algorithm. | - Limited information on specific protocols and their evaluations. |
| unequal clustering protocols [10] | - Classifies WSN clustering protocols into homogeneous and heterogeneous networks. <br> - Comparative analysis considers factors like cluster count, inter-cluster communication, CH count, clustering objects, and complexity. | - Challenges of protocols are outlined but not detailed. |
| ESO-LEACH [11] | - Explores unequal clustering techniques categorized into probabilistic, preset, and deterministic methods. <br> - Comparative analysis focuses on clustering properties and the clustering process. | - Limited information on simulation results and energy usage. |
| CFPSO [12] | - Integrates PSO to improve energy efficiency. <br> - Considers energy levels and distances for CH selection. | - High computational cost compared to LEACH. <br> - May struggle with dynamic network changes without full recalculations. |
| 5GLHN [15] | - Improves efficiency in 5G Local Home Networks using CFPSO for cell attachment. <br> - Empirical testing for resource usage and security assessment. | - Specifics of load-balancing techniques not detailed. |
| OPSO [18] | - Considers multiple criteria like remaining energy, distance, and node density. <br> - Utilizes PSO for CH selection. | - Neglects the clustering procedure, leading to energy inefficiency and lack of cluster creation acknowledgment. |

## III. PROPOSED METHODOLOGY

Introducing a groundbreaking technique called HarborSync, this research aims to enhance the capabilities of Wireless Sensor Networks (WSN). In dynamic and resource-constrained sensor network contexts, HarborSync aims to improve stability, durability, and congestion control significantly. The algorithm starts by carefully placing sensor nodes inside the target region to set the stage for future network operations.

In order to maximize network efficiency, HarborSync uses advanced methods for both initial cluster creation and continuous maintenance. The ability to delegate obligations inside clusters is a brand-new feature in HarborSync. Each node in a cluster plays a crucial role, including the leader, backup, members, and gateway. By deliberately delaying cluster head changes, HarborSync ameliorates the network's overall stability. In order to dynamically equilibrate the duties of backup nodes, old cluster heads, and new cluster heads, the system also lets in strategies for prioritizing nodes fitting into

their degree and battery life. With features like dynamic cluster reconfiguration, adaptive routing, effective aggregation, and congestion monitoring, HarborSync provides a complete answer to the myriad problems with WSNs. HarborSync coordinates the creation and upkeep of stable clusters while simultaneously negotiating and tracking congestion in real-time. The system inducts adaptive actions, such as dynamic reconfiguration [19] and efficient routing, in reaction to congestion detection to relieve network strain [20]. Finally, among the most innovative WSN algorithms, HarborSync stands out in peculiarity because it can overturn WSNs in terms of congestion control, endurance, and stability. Fig. 1 expresses all the components of the intimated algorithm stages.



Fig. 1.   Proposed methodology.

HarborSync, a urged approach, provides a novel and complex way to wangle wireless sensor networks (WSN). It begins with carefully placing sensor nodes [21] and applies a unique method to cluster formation by picking cluster heads using probability. The system has a unique approach to role assignment that takes into account vital factors, including energy, connectivity, battery life left, base station distance, sensor data quality, and priority calculations. By desegregation adaptive routing, effective aggregation techniques, dynamic cluster reconfiguration, and congestion monitoring, HarborSync offers a stiff real-time congestion control and optimization solution for dynamic and resource-constrained sensor networks.

Algorithm Components:

Initialize Network:

First, in HarborSync's operations, sensor nodes are laid strategically in the interior of the designated target zone. Every network function in the future will be ramped up in this first fundamental step. First, the N sensor node, which is the total number of sensor nodes employed by HarborSync, is consideringly positioned past the target zone.

Cluster Formation and Maintenance:

When it comes to initial cluster construction and continuing maintenance, HarborSync uses unique procedures to take care of it. This strategy aims to maximize the total number of clusters generated while also optimizing the network's overall efficiency [1-2]. Let:

- $N$ be the total number of nodes in the WSN.

- $P$ be the desired percentage of nodes that become cluster heads in each round.

- $r$ be the current round.

- $T$ be the total number of rounds.

The probability $p$ of a node becoming a cluster head in round $r$ can be determined by:

$$p = \frac{P}{1 - P \cdot (r \bmod \frac{1}{p})} \qquad (1)$$

The expected number of cluster heads $ECH$ in round $r$ is given by:

$$ECH = p \cdot N \qquad (2)$$

And the total number of clusters $C$ over $T$ rounds can be calculated as:

$$C = \sum_{R=1}^{T} ECH_r \qquad (3)$$

In this more complicated case, the predicted number of cluster heads is the total of all rounds, and each round's cluster head selection is based on probability. The total number of clusters is represented by this formula [3].

Roles within Clusters:.

HarborSync inaugurates an innovator role assignment approach, dooming roles such as Cluster Head ($CH$), Cluster

Member ($CM$), Backup Node ($BN$), and Cluster Gateway ($CG$) based on its unparalleled methodology [4].

$$CH\_Score = w1 \cdot Energy + w2 \cdot Distance\ to\ Base\ Station + w3 \cdot Connectivity + w4 \cdot Remaining\ Battery\ Life + w5 \cdot Sensor\ Data\ Quality \tag{4}$$

Where:

- Energy is the remaining energy of the node.

- Distance to Base Station: is the distance of the node to the base station.

- Connectivity is a measure of the node's connectivity in terms of communication hops.

- Remaining Battery Life is the remaining battery life of the node.

- Sensor Data Quality is the quality of the sensor data if applicable.

$w1, w2, w3, w4$, and $w5$ are weights assigned to each factor based on their relative importance. The weights should sum up to 1 to ensure normalization.

The CH selection process can then be simplified as follows:

$$Select\ CH = argmax_i\ (CH\_Score) \tag{5}$$

This means selecting the node with the highest CH score among all available nodes.

Cluster Head Changes:

To bolster overall stability, HarborSync employs a sophisticated mechanism, strategically delaying cluster head changes within the same cluster.

- $Ttransmit$ be the transmission time for a packet from a cluster head.

- $Tprocess$ be the processing time at the cluster head.

The total delay time for a Cluster Head can be expressed as:

$$DelayTime_{CH} = Ttransmit + Tprocess \tag{6}$$

Priority Calculation:

Utilizing node degree and battery life metrics, HarborSync incorporates priority calculation mechanisms. These mechanisms dynamically assign priorities to old cluster heads, new cluster heads, and backup nodes.

Formula for Priority Calculation:

$Priority = g(NodeDegree, BatteryLife, Distance)$ (7)

- $Priority$: Represents the priority assigned within the HarborSync algorithm.

- $(NodeDegree, BatteryLife, Distance$: Metrics used in the calculation within HarborSync.

- $g(NodeDegree, BatteryLife, Distance$ ): The specific function employed in HarborSync to dynamically calculate priority for nodes.

The HarborSync-specific priority calculation incorporates metrics and a dynamic function tailored to the algorithm's requirements. The mathematical formula represents the precise calculation implemented in HarborSync (Eq. (4)).

Decision Logic:

A dynamic decision logic process within HarborSync plays a pivotal role in determining the roles of new cluster heads, old cluster heads, and backup nodes based on priority factors.

Formula for Decision Logic:

$$Decision = h(Priority, OtherParameters) \tag{8}$$

Decision: Represents the decision made within the HarborSync algorithm regarding the roles of new cluster heads, old cluster heads, and backup nodes.

Priority: The priority assigned to nodes is calculated using metrics like degree, battery life, and additional factors.

Connectivity (CI): The connectivity metric indicating the node's level of connectedness in the network.

Let:

- $N$ is the total number of sensor nodes in the network.

- $L$ is the number of established links or connections between nodes.

- $R$ is the communication range of a sensor node.

- $Q$ is a measure of link quality.

A basic formula for connectivity index ($C$) can be expressed as follows:

$$CI = \frac{N(N-1)}{2} \tag{9}$$

This formula represents the ratio of the actual number of links ($L$) to the potential number of links in a fully connected network $\frac{N(N-1)}{2}$, considering undirected links). It provides a normalized measure of connectivity.

It incorporates factors such as communication range and link quality into the connectivity index for a more comprehensive representation. For instance:

$$CI = (f(x) = \sum_{i=1}^{N\infty} \sum_{j=i+1}^{N\infty} (\frac{1}{di,j}, qi,j)) \div \frac{N(N-1)}{2} \tag{10}$$

Here, $dij$ represents the distance between nodes $i$ and $j$, and $Qij$ represents the link quality between them. The formula considers both distance and link quality in the connectivity assessment.

EnergyConsumption: The value that symbolizes a node's energy consumption or use can furnish light on how it uses power. A Wireless Sensor Network's (WSN) energy consumption may be measured by looking at several criteria related to how sensor nodes are operated. E is the total energy exhausted during gearbox ($E_{tx}$), reception($E_{rx}$), and idle time ($E_{idle}$), according to a basic model:

$$E = E_{tx} + E_{rx} + E_{idle} \tag{11}$$

Here, the energy components can be defined as follows:

Energy Consumption during Transmission (*E*tx):

$$E_{tx} = P_{tx} \cdot T_{tx} \qquad (12)$$

$P_{tx}$ is the power consumption during transmission. and $T_{tx}$ is the time spent on transmission.

Energy Consumption during Reception ($E_{rx}$):

$$E_{rx} = P_{rx} \cdot T_{rx} \qquad (13)$$

$P_{rx}$ is the power consumption during reception.

$T_{rx}$ is the time spent on reception.

Energy Consumption during Idle ($E_{idle}$):

$$E_{idle} = P_{idle} \cdot T_{idle} \qquad (14)$$

$P_{idle}$ is the power consumption during the idle state.

$T_{idle}$ is the time spent in the idle state.

ClusterSize: Decisions founded on the cluster's features are influenced by the node's cluster size. Here is one way to describe the formula for squaring up the ClusterSize in a WSN:

$$\boldsymbol{ClusterSize} = \sum_{i=1}^{N} Node_i \qquad (15)$$

Here, N represents the total number of nodes in the cluster, and $Node_i$ represents the individual nodes within the cluster. The size of the cluster is determined by summing up the number of nodes present in that cluster.

Integration of HarborSync Elements:

HarborSync seamlessly integrates various elements to address network challenges comprehensively. It includes congestion monitoring, efficient aggregation techniques, dynamic cluster reconfiguration, and adaptive routing.

- Formula for Congestion Monitoring:

$$CongestionLevel = i(SensorData, Thresholds \qquad (16)$$

Here, fi is a function that calculates the congestion level for the $i-$th node based on its sensor data and predefined thresholds.

- Formula for Efficient Aggregation Techniques:

$$AggregatedData = j(SensorData, AggregationMethod) \qquad (17)$$

The function gj aggregates the sensor data for the $j^{\text{th}}$ node using a specified aggregation method.

Formula for Dynamic Cluster Reconfiguration:

$$Reconfiguration = k(ClusterStructure, CongestionLevel) \qquad (18)$$

Depending on the present cluster setup and congestion conditions, the $h_k$ function reconfigures the cluster structure dynamically.

Formula for Adaptive Routing:

$$RoutingPath = l(ClusterStructure, CongestionLevel) \qquad (19)$$

Taking into account the present cluster topology and congestion levels, the adaptive routing path for the $l-$th node is determined by the function $m_l$. In these equations, SensorData stands for sensor data, Thresholds are congestion thresholds, AggregationMethod is the data aggregation method chosen, ClusterStructure is the current cluster configuration, and CongestionLevel is the computed congestion level for a node. Taking into account the present cluster topology and congestion levels, the adaptive routing path for the $l-$th node is determined by the function $m_l$. Thresholds are predetermined levels used to determine congestion, while SensorData is the data acquired by the sensors in these formulae. The selected technique for data aggregation is AggregationMethod, the current arrangement of sensor nodes in clusters is represented by ClusterStructure, and the computed congestion level for a node is CongestionLevel.

Overall Flow:

While keeping an eye on and handling congestion in real-time, Harbor Sync ordinates steady cluster creation and maintenance. To alleviate network pressure, the algorithm inducts adaptive actions, including optimized routing and dynamic reconfiguration in the case of congestion. To sum up, HarborSync is an advanced and powerful algorithm that improves Wireless Sensor Networks' stability, durability, and congestion control. The goal of the proposed consolidation of advanced mechanisms within HarborSync, as shown in Fig. 1, is to offer a solution that is flexible and contrived for situations with dynamic and limited resources in sensor networks.

## IV. RESULT EVALUATION

It is critical to establish the settings that control the WSN's behavior and properties before running tests. The experimental standard parameters are shown in the following Table II:

TABLE II. EXPERIMENTAL SETUP

| Parameter | Values |
|---|---|
| Number of Nodes (N) | 100 - 200 |
| Area of Deployment | 500 x 500 meters |
| Simulation Time | 50 - 250 seconds |
| Pause Time for Nodes | 5 - 25 seconds |
| Max Speed of Nodes | 2 - 10 meters/second |
| Transmission Range | 25 - 250 meters |
| Transmission Power | 1 - 100 milliwatts |

*1) Stability and durability:* One way to assess stable and long-lasting clustering methods is to look at how often the cluster heads change throughout the simulation. An approach with a lower count of cluster head varies is more stable and long-lasting since extreme swings might enhance energy consumption and network overhead.

Fig. 2. Stability evaluation.

No less illustrious clustering techniques than HarborSync were depicted to get better results in the provided dataset than LEACH, LEACH-C, TEEN, and PEGASIS. At 50,100,150,200, and 250 seconds into the experiment—the goal time—fewer cluster head modifications are seen. Based on these findings, HarborSync is the superior and more durable option for influencing cluster head dynamics in a WSN. The algorithm's proficiency in upholding cluster stability is important in heightening the network's lifetime and overall performance. In Fig. 2, you can observe the outcomes.

*2) Congestion Control parameter:* Congestion control performance is a significant metric for evaluating clustering algorithms in the context of wireless sensor networks (WSNs). Data flow and WSN functioning are both negatively impinged on by congestion, which the algorithm's capacity to palm and relieve is mensurable by congestion control.



Fig. 3. Congestion control evaluation.

In Fig. 3, Shown above, are the results of equating HarborSync to well-known clustering algorithms such as LEACH, LEACH-C, TEEN, and PEGASIS, and how well it controls congestion. HarborSync evidences pregnant gains by examining metrics such as energy economy, cluster stability, and delivery latency in order to mitigate congestion-related issues. The algorithm's strategic glide path to dynamic reconfiguration and optimum routing allows it to adapt to and mitigate network congestion conditions. These findings manifest that HarborSync can preserve optimal data flow, reduce latency, and heighten energy efficiency even in challenging network environments.

*3) Energy efficiency:* Energy efficiency is An essential cadence to consider while scoping clustering algorithms for usage in WSNs. An important component in the resource-constrained WSN environment, it measures the algorithms' ability to achieve their objectives with little energy consumption. Below are the results for the energy efficiency parameters for HarborSync, LEACH, LEACH-C, TEEN, and PEGASIS at various simulation timeframes. HarborSync optimizes power utilization better than its competitors while having lower energy efficiency ratings. This acquisition shows that by carefully controlling energy resources when the network is in use, HarborSync may grow sustainability and network longevity. In the realm of WSNs, HarborSync bears out as a possible solution that achieves a fair balance between functionality and resource conservation due to its intensity on energy efficiency, which can be seen in Fig. 4 and Table III.

*4) Adaptability:* Algorithms in wireless sensor networks (WSNs) compute the ability to adjust to changing network conditions. Part of the results are the adaptability ratings for PEGASIS, TEEN, HarborSync, LEACH, and LEACH-C.



Fig. 4. Energy consumption evaluation.

TABLE III. ENERGY CONSUMPTION COMPARISON

| Time | HarborSync | LEACH | LEACH-C | TEEN | PEGASIS |
|------|-----------|-------|---------|------|---------|
| 50 | 4.35 | 7.1 | 6.08 | 5.21 | 6.84 |
| 100 | 6.83 | 13.14 | 9.75 | 8.89 | 11.22 |
| 150 | 9.74 | 16.01 | 14.38 | 13.16 | 15.79 |
| 200 | 12.88 | 19.74 | 17.12 | 18.79 | 19.31 |
| 250 | 18.75. | 22.89 | 21.74 | 20.95 | 22.66 |

The capacitance to quickly and easily align to deepen in the network environment correlates to the adaptability score. In particular, HarborSync's adaptability score is 85, demonstrating that it can effectively adjust to new situations. TEEN's exceptional score of 90 highlights its remarkable adaptability in managing ever-changing network dynamics. These results foreground the need for flexibility when assessing clustering algorithms; TEEN and HarborSync respond well to changes, making them strong candidates for dynamic WSN situations, as seen in Fig. 5.



Fig. 5. Adaptability evaluation.

*5) Packet loss rate:* In order to gauge the reliability of data transport, measures for packet loss rate are legit for wireless communication systems and wireless sensor networks (WSNs). This is the percentage of packets that flunk to reach their intended recipient an results of various algorithm can be seen in Fig. 6.

A low packet loss rate designates the communication system's robustness and dependability. Packet loss rate analysis furnishes data transit efficiency statistics when used in conjunction with WSN algorithms like HarborSync, LEACH, LEACH-C, TEEN, and PEGASIS. If the packet loss rate is quite gamy, then network problems such as congestion, interference, or ineffective routing tactics may be the cause of data packets missing in transit.



Fig. 6. Packet loss evaluation.

## V. DISCUSSION

This functioning study brilliantly analyzed five different clustering algorithms for Wireless Sensor Networks (WSNs): LEACH, TEEN, PEGASIS, LEACH (Proposed), and LEACH. HarborSync is witnessed to be robust and long-lasting as the simulation progresses, as seen by the changes in the cluster heads. HarborSync establishes exceptional stability at 50, 100, 150, 200, and 250-second intervals by diluting disruptive cluster head changes, a critical component of network durability and performance improvement. Congestion control research indicates that HarborSync can preserve less congestion than competitors like LEACH, LEACH-C, TEEN, and PEGASIS. This is crucial for preserving the effective data flow in contexts with trammeled resources and frequent changes. Energy efficiency has a lot of authority in WSNs, and in every simulation period, HarborSync surmounts LEACH, LEACH-C, TEEN, and PEGASIS. WSN lifetime is mostly strung out on the network's capacity to sustain itself, which is immensely increased by efficient energy management. Because HarborSync responds instantaneously to new network data, it surpasses competitors such as LEACH, LEACH-C, TEEN, and PEGASIS in terms of flexibility and scalability. Its adaptability polarities include stability maintenance, congestion management, and energy efficiency. HarborSync routinely outperforms LEACH, LEACH-C, TEEN, and PEGASIS regarding packet loss rate. Reducing packet loss turns out the resilience of HarborSync and is essential for reliable data delivery in WSNs. Fig. 7 demonstrates the comparative analysis using with all parameters.

Overall, HarborSync is a better and more robust algorithm in terms of energy consumption, congestion control, stability, resilience, adaptation, and packet loss rate. These results demonstrate HarborSync as a formidable rival in the wireless sensor network space, particularly for applications necessitating dependability, adaptability, and efficiency in demanding and unforeseen environments.

Fig. 7. Overall comparison evaluation.

## VI. CONCLUSION

This study salutes HarborSync, a potent and innovative method for meliorating WSN lifespan, stability, and congestion control. Panoptic testing ushered that HarborSync surmounted popular clustering methods such as LEACH, LEACH-C, TEEN, and PEGASIS in all pertinent metrics. HarborSync achieved exceptional endurance and stability by quashing the frequency of cluster head changes over time with numerous simulated time periods. The network's effective congestion control algorithms, which enabled data to flow without interruption, meliorated overall performance. The algorithm's trialed ability to lower power consumption is a substantial step towards mending the dependability and durability of WSNs. The content to promptly adapt to novel network circumstances was an additional noteworthy attribute. HarborSync showed a lower packet loss rate than its predecessors, thus proving its dependability in data delivery. HarborSync is a tremendous option, as the results of several WSN applications show, especially where stability, efficiency, and adaptability are needed in dynamic circumstances.

## VII. FUTURE SCOPE

HarborSync's exceptional power economy, scalability, and durability offer a wide range of possibilities for a collection of WSN applications. Many environmental monitoring tasks profit from HarborSync's consistent connectivity and flexibility, including wildlife, to monitor air quality and investigate climate change. Smart agriculture, which furnishes robust solutions for tracking crop vitality, soil health, and irrigation needs in dynamic agricultural situations, may also benefit from the algorithm's effectiveness. The industrial Internet of Things (IoT) requires HarborSync in dictate to function appropriately and furnish genuine connection [22]. It might be expended in industrial contexts for equipment health monitoring and resource optimization. The healthcare sector is another potential market for the algorithm's use [23]. HarborSync's stability and competence are vital for managing healthcare facilities, supervising medicine distribution, and

attesting to the dependability of patient monitoring systems. Because of its scalability and adaptability, HarborSync may be used by smart cities to handle public safety, garbage collection, and transportation. It makes cities more resilient and effective. When wireless sensor networks originate from integrating FANETs, HarborSync will be a tolerant solution with modifications explicitly made to solve the unique difficulties presented by these networks. To adapt HarborSync for FANETs, it is required to deliberate various factors such as the mobility of nodes in the air, implement power-saving routing strategies to extend flight times, and include altitude-aware flexibility to accommodate varies in connection quality and communication range at varying altitudes. Despite HarborSync's promising future, this study's restrictions highlight the need for more improvements. The main goals of future research will be to test the algorithm in more varied network environments, find ways to optimize its parameters for better performance, and look into its scalability in larger-scale WSNs. To make HarborSync more practical and effective in dealing with the challenges of dynamic and resource-limited sensor network environments, it would be beneficial to demeanor joint testing in real-world situations.

## REFERENCES

[1] S. Pathak and S. Jain, "An optimized stable clustering algorithm for mobile ad hoc networks," EURASIP J. Wirel. Commun. Netw., vol. 2017, pp. 1–11, 2017.

[2] S. Verma, S. Zeadally, S. Kaur, and A. K. Sharma, "Intelligent and secure clustering in wireless sensor network (WSN)-based intelligent transportation systems," IEEE Trans. Intell. Transp. Syst., vol. 23, no. 8, pp. 13473–13481, 2021.

[3] W. R. Heinzelman, A. Chandrakasan, and H. Balakrishnan, "Energy-efficient communication protocol for wireless microsensor networks," in Proceedings of the 33rd annual Hawaii international conference on system sciences, 2000, pp. 10-pp.

[4] W. B. Heinzelman, A. P. Chandrakasan, and H. Balakrishnan, "An application-specific protocol architecture for wireless microsensor networks," IEEE Trans. Wirel. Commun., vol. 1, no. 4, pp. 660–670, 2002.

[5] X. Liu, "A survey on clustering routing protocols in wireless sensor networks," sensors, vol. 12, no. 8, pp. 11113–11153, 2012.

[6] M. M. Afsar and M.-H. Tayarani-N, "Clustering in sensor networks: A literature survey," J. Netw. Comput. Appl., vol. 46, pp. 198–226, 2014.

[7] S. P. Singh and S. C. Sharma, "A survey on cluster based routing protocols in wireless sensor networks," Procedia Comput. Sci., vol. 45, pp. 687–695, 2015.

[8] M. T. Quasim et al., "An internet of things enabled machine learning model for Energy Theft Prevention System (ETPS) in Smart Cities," J. Cloud Comput., vol. 12, no. 1, p. 158, 2023, doi: 10.1186/s13677-023-00525-4.

[9] A. Zeb et al., "Clustering analysis in wireless sensor networks: The ambit of performance metrics and schemes taxonomy," Int. J. Distrib. Sens. Networks, vol. 12, no. 7, p. 4979142, 2016.

[10] S. Arjunan and S. Pothula, "A survey on unequal clustering protocols in wireless sensor networks," J. King Saud Univ. Inf. Sci., vol. 31, no. 3, pp. 304–317, 2019.

[11] G. K. Nigam and C. Dabas, "ESO-LEACH: PSO based energy efficient clustering in LEACH," J. King Saud Univ. Inf. Sci., vol. 33, no. 8, pp. 947–954, 2021.

[12] M. K. Hasan et al., "Constriction factor particle swarm optimization based load balancing and cell association for 5G heterogeneous networks," Comput. Commun., vol. 180, pp. 328–337, 2021.

[13] S. Bharany, S. Sharma, S. Bhatia, M. K. I. Rahmani, M. Shuaib, and S. A. Lashari, "Energy Efficient Clustering Protocol for FANETS Using

Moth Flame Optimization," Sustainability, vol. 14, no. 10, p. 6159, May 2022, doi: 10.3390/su14106159.

[14] S. Qamar et al., "Cloud data transmission based on security and improved routing through hybrid machine learning techniques," Soft Comput., pp. 1–8, 2023.

[15] M. Suresh Chinnathampy, T. Aruna, and N. Muthukumaran, "Antenna design: Micro strip patch for spectrum utilization in cognitive radio networks," Wirel. Pers. Commun., vol. 119, pp. 959–979, 2021.

[16] N. Alqahtani et al., "Deep belief networks (DBN) with IoT-based alzheimer's disease detection and classification," Appl. Sci., vol. 13, no. 13, p. 7833, 2023.

[17] M. Shuaib et al., "An Optimized, Dynamic, and Efficient Load-Balancing Framework for Resource Management in the Internet of Things (IoT) Environment," Electronics, vol. 12, no. 5, p. 1104, 2023.

[18] N. Bhalaji, "Cluster formation using fuzzy logic in wireless sensor networks," IRO J. Sustain. Wirel. Syst., vol. 3, no. 1, pp. 31–39, 2021.

[19] S. A. Devaraj, T. Aruna, N. Muthukumaran, and A. A. Roobert, "Adaptive cluster-based heuristic approach in cognitive radio networks for 5G applications," Trans. Emerg. Telecommun. Technol., vol. 33, no. 1, p. e4383, 2022.

[20] M. Suresh Chinnathampy, T. Aruna, and N. Muthukumaran, "Design and fabrication of micro strip patch antenna for cognitive radio applications," Wirel. Pers. Commun., vol. 121, pp. 1577–1592, 2021.

[21] V. K. Malhotra, H. Kaur, and M. A. Alam, "An analysis of fuzzy clustering methods," Int. J. Comput. Appl., vol. 94, no. 19, pp. 9–12, 2014.

[22] S. Alam, "Security Concerns in Smart Agriculture and Blockchain-based Solution," in 2022 OPJU International Technology Conference on Emerging Technologies for Sustainable Development (OTCON), 2023, pp. 1–6.

[23] I. Aqeel et al., "Load Balancing Using Artificial Intelligence for Cloud-Enabled Internet of Everything in Healthcare Domain," Sensors, vol. 23, no. 11, p. 5349, 2023.

# Application of Skeletal Skinned Mesh Algorithm Based on 3D Virtual Human Model in Computer Animation Design

Zhongkai Zhan*

Wuhan Vocational College of Software and Engineering, School of Arts and Media, Wuhan, 430205, China

*Abstract*—**3D virtual character animation is the core technology of games, animation, and virtual reality. To improve its visual and realistic effects, the research focused on the skeleton skinned mesh algorithm. Firstly, a three-dimensional human body model was established based on motion capture data. Then, the skin vertex weight calculation and bone skin animation design were completed for the human body model. These experiments confirm that the designed weight calculation method has a smooth weight transition and good computational stability. The designed skinned mesh algorithm outperforms its skinned mesh algorithms in accuracy, recall, and area under curve values, with a maximum area under curve value of 0.927. Its smoothness and volume retention rate are both above 90.00%, and there is no obvious collapse phenomenon. Its other objective and subjective evaluation indicators are superior to the existing advanced skinned mesh algorithms processing, and the skin effect is realistic and smooth. Overall, this study contributes to the creation of 3D virtual character animation, enhances the visual realism of virtual creation, and provides key support for the animation performance of virtual characters.**

*Keywords*—*3D virtual human; skinned mesh algorithm; weight; character animation; dual quaternion; motion capture data*

## I. INTRODUCTION

3D Virtual Character Animation (3DVCA) is a technology that uses computer to generate virtual character characters and complete animation drawing. As computer hardware, software, and graphics processing technology develop, 3DVCA plays an increasingly important role in people's work, learning, and daily entertainment life [1-2]. At present, 3DVCA has been widely applied in fields such as game development, film production, virtual reality, and augmented reality, and it belongs to a complex interdisciplinary field. Through 3DVCA, animation developers can create realistic and realistic virtual character images, enhancing user immersion and interactivity. [3-4]. However, creating highly realistic virtual human models requires rich anatomical knowledge and artistic skills, and the generation and control of virtual human actions need to consider multiple technical points such as human biomechanics, motion capture, and fusion. And bone skin binding technology still faces challenges in handling complex angle models, complex actions, and details. These unresolved technical challenges limit the application and development of 3DVCA [5]. How to build a more realistic simulation of human motion in 3D virtual character simulation has become a hot research topic. The advancement of dynamic simulation technology for 3D human models is of great significance for

the development of interdisciplinary fields such as entertainment, industry, and healthcare. And simulating human motion mainly involves the construction of dynamic models, precise motion behavior control systems, and the binding of bones and skin. In order to solve the technical difficulties in 3D virtual character simulation, the study chooses bone skin binding as the research core. Firstly, the motion capture data are analyzed to establish a 3D Virtual Human Model (3DVHM). On the basis of completing pose matching between the skin and bone models, a skin vertex weight calculation method is designed. Finally, the skeletal skin animation design is completed based on the dual quaternion mixed Skinned Mesh Algorithm (SMA). The innovation of the research is mainly reflected in the matching of skin and bone, as well as the study of skin deformation. On the one hand, the research innovatively designed a weight calculation method for skin vertices. Use a two-layer model to match and optimize poses, and calculate skin vertex weights in the area where bones affect the skin. On the other hand, research has optimized and improved the traditional linear hybrid skin Linear Blending Skinning (LBS) algorithm to avoid skin distortion and deformation, and to address the shortcomings of existing methods.

The research mainly includes four parts. Firstly, a review of the current research status of 3DVCA is completed, and summarized advanced research on 3D modeling. Then, the process of calculating skin vertex weights based on skin and bone pose matching was explained, and the steps of constructing an animation model based on the dual quaternion hybrid skin algorithm were explained; And the performance test results of the designed bone skin animation model are analyzed, and the skin effects were compared in application. Finally, the experimental results are summarized and summarized. This study is expected to contribute to the development of 3DVCA, providing a theoretical foundation and application technology for the research of computer graphics and computer animation.

## II. RELATED WORKS

3DVCA is currently a hot topic of concern in animation, film and television, and human-computer interaction. To further improve the modeling realism of virtual characters and models, scholars have conducted research on 3D modeling related technologies. The B-spline curve of a disk is an effective modeling tool based on control points, which can directly model regions or shapes with adjustable thickness.

However, the boundary curves of the shape described by the B-spline curve of the disk often exhibit self-intersection, which has a negative impact on the shape and texture of the model. Kruppa et al. designed an iterative algorithm that utilized the approximate circular skin method to handle the details of 3D model modeling. This method could generate smooth shape textures around sharp curved parts [6]. Although physical simulation can improve the detailed dynamic motion of animated characters, these are post-processing added after the overall animation modeling is completed. Wu et al. designed a new interactive framework based on position-based dynamic unified skin transformation and kinematic simulation. This framework had the advantages of controllable skin transformation and shape preservation, ensuring the efficiency, simplicity, and stability of model creation [7]. Implicit skinning in geometric interactive skinning methods is commonly used to solve skin self-collision and handle reasonable deformation at joints, requiring more user interaction to fully parameterize. Hachette et al. designed a specialized optimization framework for implicit skinning based on particle systems and dynamic shell meshes, and optimized the shape of the filled mesh to achieve reasonable skin deformation at joint rotations [8]. The same topology structure as the human body can produce the most realistic animation effects in digital clothing body animation, but its application limitations are strong. Peng et al. designed a graph convolutional network based on deep learning that could generate realistic clothing animations. This model was suitable for clothing types that did not match the body topology. Qualitative and quantitative experiments had verified that the model achieved state-of-the-art 3D clothing animation performance [9].

Mouhou et al. designed a real-time method for generating mesh deformation based on implicit skinning using spherical primitives, considering skin contact and muscle swelling simulation. This technique compensated for the shortcomings of linear hybrid skinning, allowing the model to handle network deformation reasonably while handling collisions and preserving mesh details [10]. The existing algorithms for reconstructing body surface models based on video sequences lack modeling of the internal structure of the human body. Zhao et al. created a four-dimensional animation process based on a personalized motion full anatomy digital model of videos. They estimated the internal structure of the human body using deformable human anatomy maps and simulated the movement deformation of soft tissues through smooth nonlinear transformations. These experiments confirmed that this method surpassed existing video-based body surface modeling methods [11]. The existing dynamic 3D modeling relies on data-driven learning and optimization, but has poor robustness in tracking different features in space and time. Moreover, the grid-based linear hybrid skin model has problems optimizing the network while maintaining a consistent mesh topology. Singla et al. proposed a new algorithm for reconstructing dynamic human body shapes using sparse contour information, and the robustness of this method was verified [12]. Neural networks are crucial for 3D reconstruction and new view synthesis of rigid scenes. Chen et al. designed a connection module based on neural fields, which utilized voxel-based corresponding search and

pre-computed linear hybrid skin functions. It could achieve precise correspondence between the normative space and the constituent space, effectively optimizing shape and skin weights [13]. Li D et al. designed a static skin model based on joint increment and skin weight increment to accurately predict dynamic fabric deformation. After testing in Unity game scenarios, the model achieved real-time prediction of fabric dynamics, and the network performed well in accurately capturing fine dynamic fabric deformation [14].

In summary, there have been many studies on 3D modeling, but there is still relatively little research on the improvement and optimization of bone-skin binding algorithms. Based on 3DVCA, this study conducts relevant research on the quality improvement and technical optimization of bone-skin binding algorithms.

## III. Design of Skeletal Skinned Mesh Algorithm Based on 3D Virtual Human Model

To achieve more realistic and smooth character animation effects, research is conducted on bone SMA in 3DVCA, and methods for calculating skin vertex weights and dual quaternion mixed SMA are designed.

### A. Calculation of Skin Vertex Weights Based on Skin and Bone Pose Matching

The modeling and detail processing of 3DVCA require a lot of time and human resources. In virtual character models, bone skin binding technology is responsible for connecting the character model with the bone system, ensuring that the skin can naturally follow the changes of the bones during animation [15]. However, bone skin binding technology still encounters issues such as occlusion, distortion, and skin vertex detachment when dealing with complex character models, complex actions, and details. Therefore, the study first calculates the weight of skin mesh vertices.

3DVHM includes skeletal structure, skin, muscle structure, and other details. Motion capture data are important resources for constructing and driving the motion of the 3DVHM [16]. They refer to real human motion data recorded through sensors or cameras, including joint positions, motion trajectories, and postures of the human body. Common motion capture data formats include ASF/AMC, BVH, C3D, etc. The study used ASF/AMC format to drive the 3DVHM. ASF contains skeleton information, defining the initial state of motion, while AMC is a motion data file. The separation of skeleton and motion information facilitates their matching [17]. Eq. (1) is the motion data frame in AMC format. In Eq. (1), $f_i$ represents motion data. $r_i^j$ represents the rotation information of the bone segment $j$ in frame $i$. $t_{i,x}^0, t_{i,y}^0, t_{i,z}^0$ represent the translation components of the root node in the world coordinate system.

$$\begin{cases} f_i = \left\langle t_i^0, r_i^0, r_i^1, r_i^2, ..., r_i^n \right\rangle \\ t_i^0 = \left( t_{i,x}^0, t_{i,y}^0, t_{i,z}^0 \right) \\ r_i^j = \left( t_{i,x}^j, t_{i,y}^j, t_{i,z}^j \right) \end{cases} \quad (1)$$

The study uses a skin-bone two-layer model to construct a human bone model, numbers and names the root nodes in the ASF file, and defines the hierarchical relationship between different joint points. Bones are formed between joint points, and the end joint points store the length and direction information of the bones. The root node and joint points contain 6 and 2 degrees of freedom information, respectively, to ultimately obtain a human skeleton with physical and activity characteristics. On this basis, the study uses the 3D modeling software 3D Max to complete the 3D human body modeling, and uses the Open GL function library to draw a polygonal mesh model to obtain modeling effect map.

To drive the human skin mesh model, it is necessary to establish a connection between the skin mesh and the joints or bones, that is, to calculate the weight of the influence of bones near the skin on the skin vertices. Maya software is used to establish a skin bone model that matches the skin model, and posture matching is performed with the motion bone model constructed from ASF files. The bone node information in the skin bone model is exported, and its joint node information is consistent with the definition of the bone model in the ASF file. The initial posture matching first calculates the bone length and completes the matching of limb length. Then, the angle information between different bones is corresponded one-to-one, and the joint direction information of the skin bone model is assigned to the moving bone model.

Usually, Maya software uses the nearest distance algorithm to calculate the range of action of bones on the skin. But it only considers the distance between bone and skin vertices, ignoring the hierarchical relationship of joints will increase the estimation error of the range of action. Therefore, the study extends a certain size of bone bounding boxes along the direction of the bone and perpendicular to the bone, and uses AABB bounding boxes to divide the bone model into different regions. Based on the bone bounding boxes, the impact of the bone on the skin is determined. Fig. 1 shows the skeleton bounding box.

The specific calculation of bounding boxes generally includes data pre-processing, determining expansion direction, and expanding bounding boxes. Firstly, based on the local coordinates of the joint points and bone information, the global coordinates of the joint points are obtained. The joint point coordinates of the bone are $P_i$ and $P_{i+1}$, respectively. The projection lengths of bones on different coordinate axes are represented as $x_{i+1} - x_i$, $y_{i+1} - y_i$, and $z_{i+1} - z_i$, respectively. The direction with the greatest change in length is approximately the bone's direction. The expansion coefficient of the bounding box is determined based on the human body structure. The bounding boxes of adjacent bones should intersect at a common joint point, and the expansion coefficient at the intersection is one-fifth of bone's length. The thickness of the bounding box at the bone cross-section is the thickness from the bone to the skin. However, considering the differences in bone and skin thickness in various parts of the human body, the cross-sectional expansion coefficient was set to one-third of the bone length. Fig. 2 shows the bounding box extension area.

The bone bounding box is defined as $B_i$. The criterion for determining whether the skin vertex $v_i\left(x_v, y_v, z_v\right)$ is in the bounding box is represented by Eq. (2).

$$\begin{cases} x_{\min} - \dfrac{1}{3}length \leq x_v \leq x_{\max} + \dfrac{1}{3}length \\ y_{\min} - \dfrac{1}{3}length \leq y_v \leq y_{\max} + \dfrac{1}{3}length \\ z_{\min} - \dfrac{1}{3}length \leq z_v \leq z_{\max} + \dfrac{1}{3}length \end{cases} \quad (2)$$

Fig. 3 shows the calculation principle of skin vertices. When the skin vertex is a non-joint vertex, the weight size is related to the distance between the joint and the skin vertex, and the joint weight is calculated using the linear gradient method. When the skin vertex is a vertex at a joint, it will be affected by three joints.



Fig. 1.　Schematic diagram of the skeleton bounding box.

Fig. 2.    Schematic diagram of the bounding box extension area.



(a)Joint                                              (b)Non-joint

Fig. 3.    Calculation principle of skin apex.

If the skin vertices are non-joint vertices, the weight calculation is represented by Eq. (3). $w_a$ and $w_b$ respectively represent the weights of the joint points at both ends of the skeleton within the bounding box. $D_1$ and $D_2$ respectively represent the projection of the distance from the skin vertex to the joint points at both ends in the bone direction.

$$\begin{cases} \dfrac{w_a}{w_b} = \dfrac{D_2}{D_1} \\ w_a = \dfrac{D_2}{D_1 + D_2} \\ w_b = \dfrac{D_1}{D_1 + D_2} \end{cases} \tag{3}$$

If the skin vertex is a vertex at the joint, the weight calculation is represented by Eq. (4). $w_a, w_b, w_c$ are the joint weights.

$$\begin{cases} w_a : w_b : w_c = \dfrac{1}{D_1} : \dfrac{1}{D_2} : \dfrac{1}{D_3} \\ w_a = \dfrac{D_2 D_3}{D_1 D_2 + D_2 D_3 + D_1 D_3} \\ w_b = \dfrac{D_1 D_3}{D_1 D_2 + D_2 D_3 + D_1 D_3} \\ w_c = \dfrac{D_1 D_2}{D_1 D_2 + D_2 D_3 + D_1 D_3} \end{cases} \tag{4}$$

### B. Animation Model Design Based on Dual Quaternion Hybrid Skinned Mesh Algorithm

Skin deformation is the process of transforming the skin of one object onto another object, commonly used in applications such as character animation, model fusion, and deformation [18]. It is very important in 3DVCA and plays a crucial role in achieving realistic and smooth character animation effects. SMA is an important algorithm in skin deformation

technology. It associates the skeletal system of the 3D human model with the skin mesh, updates the position of the skin mesh vertices in the world coordinate system, and enables the character's skin to naturally follow the movement changes of the bones during animation [19]. The performance of SMA is related to the deformation and transformation of animated characters, determines the adjustment of model details and quality, and directly affects the efficiency and performance of rendering. The selection of SMA is very important.

Homogeneous coordinates and matrix forms are used to describe the coordinate changes in three-dimensional space. The homogeneous coordinate change is represented by Eq. (5). $(x, y, z)$ is the original coordinate. $(x', y', z')$ is a homogeneous coordinate. $h$ means the scaling factor.

$$(x, y, z) \Rightarrow \begin{cases} x = x'/h \\ y = y'/h \\ z = z'/h \end{cases} \tag{5}$$

The skin vertex undergoes translation or rotation transformation, and the transformed skin vertex $v'$ is represented by Eq. (6). $\mathbf{M}$, $\mathbf{T}$, and $\mathbf{R}$ are skin vertex change, translation transformation, and rotation transformation matrices, respectively. The rotation matrix can be decomposed into the rotational changes of a point around three coordinate axes.

$$\begin{cases} v' = \mathbf{M} \cdot v \\ \mathbf{M} = \mathbf{T} \cdot \mathbf{R} \end{cases} \tag{6}$$

The motion of the designed 3D virtual is based on the global coordinate system, and the basic position of the model is determined by the root nodes and joint points of the skeleton. The motion of the joint points is defined relative to the position of the parent node, and these two affect each other. So the base coordinate system and local coordinate system in Fig. 4 are introduced to describe this bone model. The base coordinate system refers to the coordinate system established with the sacral joint point as the origin in the initial pose of a three-dimensional virtual human. The local coordinate system is the rotational information defined relative to the global coordinate system.

The change in the coordinate line of any joint point is represented by Eq. (7). $\varphi$, $\alpha$, and $\beta$ are the angles at which a node rotates around $z, x, y$ axes. $G_i$ represents the rotation offset matrix of the related node relative to the parent node in the world coordinate system.

$$G_i = R_i^z(\varphi) \cdot R_i^x(\alpha) \cdot R_i^y(\beta) \cdot T_i \cdot \left[R_i^z(\varphi)\right]^{-1} \cdot \left[R_i^x(\alpha)\right]^{-1} \cdot \left[R_i^y(\beta)\right]^{-1} \tag{7}$$

The global change matrix $M_i$ is represented by Eq. (8). $f(i)$ is the parent node number of the joint point. $F$ means a set of all node numbers. $G_0$ represents the global affine transformation matrix of the root node, $G_0 = R_0 \cdot T_0 \cdot R_0^{-1}$.

$$M_i = G_0 \cdot \ldots \cdot G_{f(i)} \cdot G_i = \prod_{k=0, k \in F}^{i} G_k \tag{8}$$

SMA uses transformation matrices to bind skin and bones, and the calculation of LBS is represented by Eq. (9). $v'$ and $v_0$ are the positions of skin vertices after and before the transformation, respectively. $n$ represents the number of bones that affect the current skin. $D_i$ means the global affine transformation matrix from the local coordinates of the skeleton to the global coordinates in the initial state of the model, and its calculation is similar to $M_i$.

$$v' = \sum_{i=1}^{n} w_i M_i D_i^{-1} v_0 \tag{9}$$

The transformation process from skin vertices in the global coordinate system to the local coordinate system in the initial state is the inverse transformation of Eq. (9). After transformation, AMC data are used for motion transformation to complete skin deformation. The deformation calculation is consistent with Eq. (8), and the values of the rotation offset matrix are analyzed and read based on AMC data. Fig. 5 shows the working mechanism of the entire LBS.



Fig. 4. Global coordinate system and local coordinate representation.

Fig. 5. Schematic diagram of skin deformation.

LBS is widely used in skin deformation calculation, but there are still significant application defects in LBS. LBS calculates bone influence weights through linear interpolation. When vertices are affected by multiple bones or overlap with multiple bones, it is easy to cause volume scaling or cross deformation problems [20]. Therefore, the study improves on the basis of LBS and designs a dual quaternion hybrid SMA using dual quaternions. The dual quaternion $\hat{q}$ consists of pairs of even numbers and quaternions, the two expressions for even numbers are given in Eq. (10). One way of expressing this can be understood as a quaternion in which all elements of the formula are pairs of even numbers, and the other expression can be understood as a quaternion in which all elements are pairs of even numbers. $w, \hat{x}, \hat{y}, \hat{z}$ are even numbers. $q, q_\varepsilon$ are quaternions, representing the real and dual parts respectively, the real part contains the rotation information. $\varepsilon$ means dual units. $i, j, k$ are orthogonal unit vectors.

$$\begin{cases} \hat{q} = w + i\hat{x} + j\hat{y} + k\hat{z} \\ \hat{q} = q + \varepsilon q_\varepsilon \end{cases} \tag{10}$$

The unit dual quaternion of the three-dimensional space vector $\vec{r} = (r_0, r_1, r_2)$ is represented by Eq. (11). When $q_\varepsilon = 0$, $\hat{q}\hat{r}\hat{q}^*$ means the rotational transformation of the rigid body. $\hat{q}^*$ and $\overline{\hat{q}}$ are conjugate and dual conjugate of dual quaternions, respectively.

$$\hat{r} = 1 + \varepsilon \left( r_0 i + r_1 j + r_2 k \right) \tag{11}$$

The unit dual quaternion of the translation transformation vector $\vec{t} = (t_0, t_1, t_2)$ of a rigid body is represented by Eq. (12). $\hat{t}\hat{v}\overline{\hat{t}^*}$ is the translation transformation of a rigid body.

$$\hat{t} = 1 + \frac{\varepsilon}{2} \left( t_0 i + t_1 j + t_2 k \right) \tag{12}$$

The process of a rigid body rotating first and then translating in three-dimensional space is represented by Eq. (13), which means that the dual quaternion $\hat{t}\hat{q}$ can achieve the process of rotating first and then translating.

$$\hat{t}\left( \hat{q}\hat{v}\overline{q^*} \right)\overline{\hat{t}^*} = \left( \hat{t}\hat{q} \right)\hat{v}\overline{\left( \hat{t}\hat{q} \right)^*} \tag{13}$$

The transformation matrix in LBS is converted into a dual quaternion, and the translation information in the matrix is defined as $a = (a_{14}, a_{24}, a_{34})$. The transformation between the translation matrix and the unit dual quaternion is represented by Eq. (14).

$$(a_{14}, a_{24}, a_{34}) = (t_0, t_1, t_2) \tag{14}$$

The coordinate of the skin vertex position $p'$ after linear mixing of dual quaternions is represented by Eq. (15). $p$ is the position coordinate of the current skin mesh vertex. $w_i$ is

the weight. Eq. (15) represents the process of mixing and unitizing the global transformation matrix according to the weight values after transforming it into dyadic quaternions.

$$p^{'} = \frac{\sum_{i=1}^{n} w_i \hat{q}_i}{\left\| \sum_{i=1}^{n} w_i \hat{q}_i \right\|} p \qquad (15)$$

Through this change, LBS converts the linear transformation of the global transformation matrix $E_i = M_i D_i^{-1}$ into linear mixing of dual quaternions. The mixing of dual quaternions was completed based on the weight size. The mixed dual quaternions are unitized and then transformed into a transformation matrix, which are multiplied by the coordinates of the skin vertices to update the skin vertices. Fig. 6 shows the workflow of the entire dual quaternion mixed SMA.



Fig. 6. Work flow of dual quaternion mixed skinned mesh algorithm.

## IV. PERFORMANCE TESTING AND APPLICATION ANALYSIS OF COMPUTER SKELETAL SKINNED MESH ALGORITHM ANIMATION DESIGN

To verify the application effect of the bone SMA designed in the design of 3DVHMs, relevant performance tests and effect evaluation experiments were designed, and the results were analyzed and discussed.

### A. Performance Testing of Weight Calculation Method Based on Posture Matching and Dual Quaternion Hybrid Skinned Mesh Algorithm

Experimental Environment Setting: The experiment was conducted on an operating system running on Windows 7, i7 CPU, GTX1060 GPU, and 16GB of memory. All algorithm models were implemented using C++and Python programming. Firstly, the designed posture matching-based Extended Bounding Box (EBB) weight calculation method were compared with other weight calculation methods for skin vertices, including Heat Balance Principle Algorithm (HBPA), Approximation Calculation of Skeletal Projection (ACKP), and Hand brush weight (HBW). Weight smoothness and stability were selected as the evaluation indicators for testing. Fig. 7 shows the experimental results. The EBB-based weight calculation method, which corresponded to the pose of the sports skeleton and skin skeleton models, performed well in the smoothness of skin surface weight. After 180 iterations, its weight smoothness finally stabilized at 90.08%, which was significantly better than other methods and 20.05 percentage points higher than HBPA of 70.03%. Meanwhile, this method had better stability in weight calculation and lower error rates for different models.

To conduct performance testing on SMA, a program was written using C++ language in the VC++6.0 software to implement the required algorithms for testing. The designed Dual Quaternion Blending Skinning (DQBS) was compared with Spherical Hybrid Skinning Deformation Algorithm (SHSD), Skinning Algorithm Based on Position-Based Dynamics (SA-PBD), and traditional linear hybrid SMA. 1686 sets of 3D virtual character models created by the design and modeling team of a certain animation development company were selected as the test dataset, which was divided into a training set and a testing set in a 7:3 ratio.



Fig. 7. Comparison of weight smoothness and stability of different weight calculation methods.

Comparing the precision and recall of four SMA algorithms, Fig. 8 shows the PR curves of different algorithms. PR is a curve drawn with precision as the vertical axis and recall as the horizontal axis. When the precision of the designed DQBS was 80% and 90%, the corresponding recall was 75.89% and 85.98%, respectively. In contrast, when the precision was 90, the recall of SHSD, SA-PBD, and LBS was 61.98%, 57.74%, and 50.16%, respectively. The precision and recall of SHSD, SA-PBD, and LBS were significantly lower than those of the study design method.



Fig. 8. Comparison of accuracy and recall rates of different skin algorithms.

The Receiver Operating Characteristic (ROC) curve is selected to evaluate different SMAs, and Fig. 9 shows the results. AUC is the area below ROC, used to measure the stability of model performance and overall effectiveness. Different SMAs had different AUC values, and the designed DQBS had the highest AUC, reaching 0.927. SA-PBD had the smallest AUC, only 0.634. Overall, the designed SMA performs better.



Fig. 9. Comparison of ROC curves for different skin algorithms.

### B. Evaluation of Skinned Effects for 3d Virtual Human Based on Dual Quaternion Hybrid Skin

The simulation system for skin animation uses Visual Studio 2010 as the development platform and motion capture data ASF/AMC files designed and developed by Acclaim Games. The skin effects of SMA with inherent defects were compared, and the knee bone bending, armpit bending, and elbow bending with larger joint rotation were selected. From an objective evaluation perspective, smoothness and volume retention rate were selected as evaluation indicators for effect comparison. Fig. 10 shows the experimental results. After 50 iterations, the smoothness and volume retention of the SMA in three different parts were all above 90.00%. The designed model did not show obvious collapse, and the volume loss of the model was relatively small.



(a) Smoothness



(b) Volume retention

Fig. 10. Comparison of smoothness and volume retention of different skinning algorithms.

The study selected objective and subjective evaluation indicators, including model generation time, bone generation time, skin time, memory usage, stability, texture preservation ability, usability, and scalability. Table I shows the statistical results. The calculation time of DBQS was short, and the skinning time was only 0.49S. Its entire weight allocation, vertex transformation, and interpolation processes were calculated at a fast speed, and the memory usage was only 1247K, indicating good performance. In addition, the stability and texture retention ability of DBQS reached 90.63% and 80.68%, respectively, significantly higher than other SMAs. The subjective evaluation indicators had obvious advantages in usability and scalability, both exceeding 85%.

Finally, animation fluency and appearance realism were selected to rate the 3DVHM. A total of 30 users were gathered to evaluate the generation models of different SMAs from the perspectives of arms, legs, abdomen, and head. Fig. 11 shows the experimental results. Authenticity is the evaluation of whether the final effect of a model is realistic and whether it can provide users with an immersive feeling. Fluency is the evaluation of whether the animation effect is smooth and natural, especially in complex deformations and actions, whether it can maintain good continuity and fluency. In Fig. 11, the fluency and authenticity scores of DBQS were relatively high in different parts, with median scores ranging from 70 to 80. Its overall effect is good.

TABLE I.    COMPARISON OF SUBJECTIVE AND OBJECTIVE EVALUATION OF DIFFERENT SKIN ALGORITHMS

| Evaluating indicator | DBQS | SA-PBD | SHSD | LBS |
|---|---|---|---|---|
| Model generation time | 13.35S | 15.19S | 12.73S | 18.25S |
| Bone formation time | 1.20S | 1.32S | 1.29S | 1.41S |
| Skin time | 0.49S | 0.52S | 0.74S | 0.96S |
| Internal memory | 1247K | 2871K | 3217K | 4126K |
| Stability | 90.63% | 63.45% | 73.01% | 72.35% |
| Texture retention | 80.68% | 68.97% | 76.17% | 62.06% |
| Expandability | 94.73% | 72.01% | 69.39% | 84.14% |
| Ease of use | 86.21% | 74.58% | 66.37% | 77.69% |



(a) Animation fluency



(b) Appearance authenticity

Fig. 11.  Animation fluency and appearance authenticity scores of different skin algorithms.

## V. CONCLUSION

The 3D virtual human skeleton skin animation technology is the key to 3DVCA creation. To achieve more realistic 3D virtual human animation effects and achieve precise bone skin deformation, a skin vertex weight calculation and hybrid SMA for virtual characters were studied and designed. These experiments confirm that the designed EBB-based weight calculation algorithm performs well in weight smoothness and stability, with the highest smoothness reaching 90.08%, which is 20.05 percentage points higher than the lowest value. After 50 iterations, the designed SMA has higher smoothness and volume retention in different parts, and there is no obvious collapse phenomenon and the volume loss rate is relatively small. The calculation time of DBQS is short, the skinning time is only 0.49S, and the memory usage is only 1247K. Its stability and texture retention ability reaches 90.63% and 80.68%, respectively. Its subjective evaluation indicators have obvious advantages in usability and scalability. DBQS has high scores for fluency and authenticity in four different parts, with median scores ranging from 70 to 80. Its overall effect is more realistic. The designed SMA has obvious advantages in the application of 3D virtual human animation design. The skin of virtual characters can naturally follow the movement changes of bones during animation. This is a more realistic presentation of character animation, which is the technical key to improving animation realism, smoothness, and rendering efficiency. However, research has adopted a double-layer structure to model virtual humans, and future research can adopt more accurate skeleton extraction algorithms. The study did not involve changes in the muscles of the characters, nor did it simulate facial expressions and expressions, which can be a focus of future research work.

## VI. SYMBOL INDEX

All conforming explanations in the text are shown in Table II.

TABLE II. INTERPRETATION OF SYMBOLS USED IN THE STUDY

| Notation | Interpretations |
|---|---|
| 3DVCA | 3D Virtual Character Animation |
| 3DVHM | 3D Virtual Human Model |
| SMA | Skinned Mesh Algorithm |
| EBB | Extended Bounding Box |
| HBPA | Heat Balance Principle Algorithm |
| HBW | Hand brush weight |
| DQBS | Dual Quaternion Blending Skinning |
| SA-PBD | Skinning Algorithm Based on Position-Based Dynamics |
| SHSD | Spherical Hybrid Skinning Deformation Algorithm |
| ROC | Receiver Operating Characteristic |
| LBS | Linear Blending Skinning |

## REFERENCES

[1] Demirel H O, Salman A, Vincent G D. Duffy. Digital human modeling: A review and reappraisal of origins, present, and expected future methods for representing humans computationally. International Journal of Human–Computer Interaction, 2022, 38(10): 897-937.

[2] Yang Y. The skinning in character animation: A survey. Academic Journal of Computing & Information Science, 2022, 5(4): 4-17.

[3] Zhang J. Survey of Skinning Method in 3D Character Animation. Academic Journal of Computing & Information Science, 2023, 6(9): 110-114.

[4] Amin, S. N., Shivakumara, P., Jun, T. X., Chong, Y. K., Zan, D. L. L., Rahavendra, R. An Augmented Reality-Based Approach for Designing Interactive Food Menu of Restaurant Using Android.Artificial Intelligence and Applications. 2023, 1(1): 26-34.

[5] Liu C, Wang A, Bu C, Wang W, Fang Z, He S. Reconstructing detailed human body using a viewpoint auto-changing RGB-D sensor for rescue operation. IEEE Sensors Journal, 2022, 22(13): 13262-13272.

[6] Kruppa K, Kunkli R, Hoffmann M. A skinning technique for modeling artistic disk B-spline shapes. Computers & Graphics, 2023, 115(6): 96-106.

[7] Wu Y, Umetani N. Two-Way Coupling of Skinning Transformations and Position Based Dynamics. Proceedings of the ACM on Computer Graphics and Interactive Techniques, 2023, 6(3): 1-18.

[8] Hachette O, Canezin F, Vaillant R, Mellado N, Barthe L. Automatic shape adjustment at joints for the implicit skinning. Computers & Graphics, 2022, 102(1): 300-308.

[9] Peng T, Kuang J, Liang J, Hu X, Miao J, Zhu P, Li L. GSNet: Generating 3D garment animation via graph skinning network. Graphical Models, 2023, 129(5): 101197-101206.

[10] Mouhou A A, Saaidi A, Yakhlef M B, Abbad K. Virtual hand skinning using volumetric shape. International Journal of Computer Aided Engineering and Technology, 2023, 18(3): 77-96.

[11] Zhao R, Mi J, Jiang Y, Chen Z, Wang H. 4DFA: Four-Dimensional Full-Anatomy Reconstruction of Individualized Digital Human Models Based on Motion Videos. Advances in Engineering Technology Research, 2023, 4(1): 269-269.

[12] Singla K, Nand P. Reconstructing dynamic human shapes from sparse silhouettes via latent space optimization of Parametric shape models. Turkish Journal of Electrical Engineering and Computer Sciences, 2023, 31(2): 295-311.

[13] Chen X, Jiang T, Song J, Rietmann M, Geiger A, Black M J, Hilliges O. Fast-SNARF: A fast deformer for articulated neural fields. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023, 45(10): 11796 – 11809.

[14] Li Y D, Tang M, Chen X R. D-Cloth: Skinning-based Cloth Dynamic Prediction with a Three-stage Network. COMPUTER GRAPHICS forum. 2023, 42(7): 14937-14939.

[15] Santesteban I, Otaduy M, Thuerey N, Casas D. Ulnef: Untangled layered neural fields for mix-and-match virtual try-on. Advances in Neural Information Processing Systems, 2022, 35(5): 12110-12125.

[16] Wu Y, Chen Z, Liu S, Ren Z, Wang S. Casa: Category-agnostic skeletal animal reconstruction. Advances in Neural Information Processing Systems, 2022, 35(7): 28559-28574.

[17] Lu Y, Yu H, Ni W, Song L. 3D real-time human reconstruction with a single RGBD camera. Applied Intelligence, 2023, 53(8): 8735-8745.

[18] Zheng Z, Yu T, Liu Y, Dai Q. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. IEEE transactions on pattern analysis and machine intelligence, 2021, 44(6): 3170-3184.

[19] Chandran P, Zoss G, Gross M, Gotardo P, Bradley D. Shape Transformers: Topology-Independent 3D Shape Models Using Transformers.Computer Graphics Forum. 2022, 41(2): 195-207.

[20] Pham J, Wyetzner S, Pfaller M R, Parker D W, James D L, Marsden A L. svMorph: Interactive geometry-editing tools for virtual patient-specific vascular anatomies. Journal of Biomechanical Engineering, 2023, 145(3): 31001-31008.

# Applying Computer Vision and Machine Learning Techniques in STEM-Education Self-Study

Rustam Abdrakhmanov[1], Assyl Tuimebayev[2], Botagoz Zhussipbek[3], Kalmurat Utebayev[4],
Venera Nakhipova[5], Oichagul Alchinbayeva[6], Gulfairuz Makhanova[7], Olzhas Kazhybayev[8]

International University of Tourism and Hospitality, Turkistan, Kazakhstan[1]
Boston University, Boston, USA[2]
Korkyt Ata Kyzylorda University, Kyzylorda, Kazakhstan[3, 7]
M. Auezov South Kazakhstan University, Shymkent, Kazakhstan[4, 6]
South Kazakhstan Pedagogical University named after O. Zhanibekov, Shymkent, Kazakhstan[5]
Astana IT University, Astana, Kazakhstan[8]

*Abstract*—In this innovative exploration, "Applying Computer Vision Techniques in STEM-Education Self-Study," the research delves into the transformative intersection of advanced computer vision (CV) technologies and self-directed learning within Science, Technology, Engineering, and Mathematics (STEM) education. Challenging traditional educational paradigms, this study posits that sophisticated CV algorithms, when judiciously integrated with modern educational frameworks, can profoundly augment the efficacy of self-study models for students navigating the increasingly intricate STEM curricula. By leveraging state-of-the-art facial recognition, object detection, and pattern analysis, the study underscores how CV can monitor, analyze, and thereby enhance students' engagement and interaction with digital content, a pioneering stride that addresses the prevalent disconnect between static study materials and the dynamic nature of learner engagement. Furthermore, the research illuminates the critical role of CV in generating personalized study roadmaps, effectively responding to individual learner's behavioral patterns and cognitive absorption rhythms, identified through meticulous analysis of captured visual data, thereby transcending the one-size-fits-all educational approach. Through rigorous qualitative and quantitative research methods, the paper offers groundbreaking insights into students' study habits, proclivities, and the nuanced obstacles they face, facilitating the creation of responsive, adaptive, and deeply personalized learning experiences. Conclusively, this research serves as a clarion call to educators, technologists, and policy-makers, emphatically demonstrating that the thoughtful application of computer vision techniques not only catalyzes a more engaging self-study landscape but also holds the latent potential to revolutionize the holistic STEM education ecosystem.

*Keywords—Load balancing; machine learning; server; classification; software*

## I. INTRODUCTION

In the rapidly evolving educational landscape, traditional teaching methodologies are incessantly being re-evaluated and challenged, particularly in Science, Technology, Engineering, and Mathematics (STEM) disciplines. The advent of digital technology has reshaped pedagogical strategies, heralding new approaches like self-directed learning, which has gained prominence for fostering students' autonomy and responsibility in the learning process [1]. However, maximizing the efficacy of self-study in STEM education requires addressing intrinsic complexities and diverse student engagement methodologies [2]. This research aims to bridge this gap by harnessing computer vision (CV) techniques, offering a transformative perspective on enhancing self-study's effectiveness in STEM education.

STEM fields, inherently multifaceted and dynamic, demand educational approaches that not only convey complex concepts but also adapt to individual cognitive styles and paces [3]. Traditional self-study, while offering flexibility, often falls short of this adaptability, leading to learner frustration and sub-optimal learning outcomes [4]. Computer vision's potential in education, particularly in monitoring and responding to student engagement and facilitating personalized learning trajectories, remains largely underexplored [5].

Computer vision, a field that grants computers a high-level understanding of digital images and videos, is traditionally aligned with applications in security, surveillance, and detection [6-7]. However, its implications extend profoundly into educational realms. Through detailed visual data analysis, CV holds the promise of decoding student engagement patterns, providing educators with nuanced insights into the often imperceptible aspects of self-study behaviors that either catalyze or hinder learning. This research pivots around the innovative application of CV in capturing and analyzing these intricate behavioral nuances, thereby guiding the development of more responsive and adaptive self-directed learning models.

Integrating CV into education, especially within STEM, poses unique challenges and opportunities. The precision required in STEM subjects translates to the necessity for educational resources to adapt in real-time to students' understanding, ensuring concepts are neither misinterpreted nor oversimplified [8]. By employing CV techniques, such as facial recognition and object detection, it becomes feasible to analyze students' interaction with educational content, thereby tailoring materials and study paths that resonate with individual learning approaches, an advancement far beyond the capabilities of traditional educational software [9].

Moreover, the role of CV in tracking and analyzing engagement brings a new dimension to educators' understanding of student performance. Conventional assessment methods offer only summative feedback, often

neglecting the formative aspects of a learner's journey [10]. Through continuous and non-intrusive monitoring, CV provides a wealth of formative feedback, empowering educators to make informed, timely interventions and students to gain awareness of their learning habits.

Significantly, the ethical considerations of utilizing CV in education are paramount, entailing careful navigation. Privacy concerns, data security, and the consent of the involved parties are crucial factors that educators and technologists must prioritize. Establishing robust ethical protocols and transparent operational guidelines ensures the responsible use of CV in educational settings, safeguarding participants' fundamental rights while harnessing technology's benefits [11].

In light of the above, this study ventures into an interdisciplinary examination of how CV can revolutionize self-study within STEM education. It builds on existing literature that outlines the theoretical frameworks of self-directed learning and delves into empirical evidence supporting the integration of advanced technologies in education [12]. By establishing a symbiotic relationship between CV technology and educational pedagogy, this research underscores a forward-thinking approach to cultivating STEM competencies, proposing a model that respects individual cognitive differences and celebrates personalized educational journeys.

Through this paper, we invite educators, technologists, and policy-makers to envision a future where technology and education converge to offer enriched, student-centric learning experiences. By navigating the technical, pedagogical, and ethical terrains of this integration, we aim to construct a comprehensive understanding that could fundamentally transform the way STEM education perceives and leverages self-study.

## II. RELATED WORKS

The synthesis of technology and education, especially in self-directed STEM learning, has instigated a plethora of research, with various studies corroborating the transformative potential of integrating advanced technological frameworks, such as computer vision (CV), into educational models. These scholarly pursuits, encompassing a diverse range of insights and findings, lay the groundwork for understanding the trajectory and implications of utilizing CV in self-regulated learning environments.

Starting with the broader impacts of technology in education, studies have indicated a paradigm shift in instructional strategies, emphasizing the need for more learner-centric approaches facilitated by technology [13]. In the context of STEM education, researchers have highlighted the necessity for innovative methods that cater to enhancing students' critical thinking and problem-solving skills, proposing that digital technologies can bridge the gap between theoretical knowledge and practical application [14].

The concept of adaptive learning, pivotal to this discussion, leverages technology to tailor educational experiences to individual needs. One study [15] provide insights into adaptive learning systems' role in promoting cognitive growth, arguing that these systems accommodate diverse learners' profiles, thereby fostering a more inclusive learning environment.

However, the challenge remains in effectively tracking and interpreting individual learner interactions and responses in real-time, a gap that computer vision promises to address [16].

Computer vision's foray into educational strategies marks a relatively new venture. Its application has been predominantly explored in surveillance, recognition systems, and user interaction tracking in various sectors [17]. Within educational research, studies have often circumscribed their focus to online learning environments, utilizing simple CV techniques for user log-in and basic interaction [18]. However, more nuanced applications of CV, such as emotion recognition, behavioral analysis, and engagement tracking, are emerging themes in contemporary literature [19].

Next study [20] delve into the potential of CV in recognizing and interpreting human emotions, an aspect crucial for personalizing learning experiences. They argue that emotional states play a significant role in learning efficiency, with certain emotional conditions favoring the absorption and retention of new information. Incorporating CV into educational platforms could thus allow for real-time adaptation based on learners' emotional cues, providing immediate feedback or altering content presentation to enhance learning efficacy [21].

Furthering the discourse on personalized learning, researchers have explored data-driven approaches. For example, [22] highlights the importance of learning analytics in understanding students' learning processes. They discuss how data obtained from students' digital footprints on learning platforms can inform more personalized and effective teaching strategies. This data-centric approach aligns with the capabilities of CV in capturing and analyzing extensive datasets of learner interactions and behaviors [23].

In the realm of self-directed learning, especially in online and digital contexts, maintaining student engagement and motivation is paramount. Studies by Chen, Lambert, and Guidry [24] underscore the challenge educators face in keeping students engaged with digital platforms. CV's potential to monitor visual cues and physical responses presents unprecedented opportunities for understanding and enhancing student engagement on a much finer, more personalized scale [25].

The integration of CV in education also extends to practical skill-based learning in STEM. For instance, research on laboratory learning indicates that CV can significantly enhance remote laboratory experiences, a critical component of STEM education. Gravier, Fayolle, Bayard, Ates, and Lardon [26] have explored these prospects, emphasizing that CV can facilitate more interactive and hands-on experiences in a virtual environment.

Despite the promising advancements, the literature consistently echoes the ethical implications of employing CV in educational settings. Privacy concerns, particularly with facial recognition and behavioral tracking, are prevalent [27]. Next research [28] stresses the need for robust privacy protection frameworks, emphasizing informed consent, data security, and transparency in how monitoring technologies are used in education. These considerations are crucial in ensuring

the ethical integrity of integrating any form of surveillance or tracking technology into learning environments.

Conclusively, the body of work surrounding the integration of computer vision in education outlines a landscape ripe with potential yet requires careful navigation concerning ethical, technical, and pedagogical constraints. This research contributes to this ongoing scholarly dialogue, contextualizing the application of CV within the specific challenges and opportunities presented by self-directed STEM education [29-32]. Through an interdisciplinary lens, this study seeks to build upon the foundations laid by existing literature [33-37], proposing an innovative convergence of CV technology and educational pedagogy to enhance the quality and effectiveness of self-study in STEM disciplines.

## III.    MATERIALS AND METHODS

In order to investigate the central research query, several "STEM Workshops: Python and Raspberry Pi Practical Activity" were organized as a precursor to the main experimental procedure. These preliminary sessions were instrumental in gathering the necessary data for the creation and subsequent validation of the RASEDS, directly contributing to the resolution of the initial research query. Once the efficacy of RASEDS was confirmed, the data derived from the system were harnessed to develop a predictive model for student performance in STEM subjects. Subsequently, this predictive mechanism was designed to suggest customizable learning resources, tailored to forecasted performance trends.

Our research adopted a quasi-experimental design to ascertain whether introducing personalized educational resources, recommended through RASEDS, could significantly improve student involvement and confidence in STEM-related tasks, thereby providing comprehensive answers to the second and third research inquiries. The sequence of research activities is graphically represented in Fig. 1.

To meticulously record the nuances of each learner's practical engagement, we strategically installed cameras to film their hands-on interaction with the educational materials, a critical component for the RASEDS's engagement detection mechanism. Care was taken in choosing camera perspectives that would clearly record the learners' hands and the instructional tools they used. Mindful of the ethical considerations when filming individuals, especially those underage, we established rigorous measures to secure informed consent from all attendees or their legal guardians (for those younger than 18. This measure was pivotal in maintaining ethical standards concerning the visual content that included identifiable participant imagery.

In the aftermath of these sessions, we collected 4,515 photographs. These were methodically divided into primary and secondary datasets, following an 80:20 split. Consequently, we allocated 3,612 photographs for initial training purposes and reserved 903 as a subsequent test collection. These images are integral to the training phase of the YOLOR model, representing a significant stride towards achieving our research's overarching goals.



Fig. 1. Flowchart of the proposed system.

Following the workshop's end, the participants showcased their STEM initiatives, firmly rooted in the Internet of Things (IoT) sphere. Utilizing the insights gained about sensors and coding principles throughout the workshop, the students embarked on devising creative approaches to tangible issues through the application of IoT. Their ventures spanned various concepts, from intelligent domestic setups and energy-saving configurations to automated methods promoting greener lifestyles and operational spaces.

Upon the conclusion of the project presentations, each was subjected to a thorough analysis conducted by two connoisseurs within the STEM domain. The assessment protocol was grounded in the principles specified by the Creative Product Analysis Matrix (CPAM) approach, involving three broad categories and nine evaluative markers, elaborately itemized in Table I. This technique of appraisal, validated in its efficacy by Besemer in 1998, guaranteed an exhaustive and precise examination of the students' endeavors.

TABLE I.        DIALOG TESTS IN CLASSROOM

| Scale | Indicator |
|---|---|
| Novelty | Original |
| | Amazing |
| Resolution | Valuable |
| | Useful |
| | Understandable |
| Synthesis | Organic |
| | Elegant |
| | Good |

Evaluation was carried out utilizing a five-point Likert scale, enabling a detailed interpretation of each project's merits and areas for improvement. The consistency of scoring between the two experts was confirmed, with a correlation coefficient marking between 0.68 and 0.84. This high degree of concordance underscored the substantial agreement in their assessments, bolstering the integrity of the evaluation phase. Such a metric reinforced the consistency and trustworthiness of the ratings given, laying a dependable groundwork for the authentic data essential for substantiating the predictive model of STEM learning outcomes.

In response to the intricate and ever-evolving facets of STEM activity-based learning, we pioneered a system known as the Real-time Automated STEM Engagement Detection System. This system is designed to autonomously and instantaneously gauge students' engagement levels. At its core, RASEDS utilizes cutting-edge object detection, particularly the YOLOR method, to pinpoint the presence of students' hands and all associated educational materials engaged during the tasks. This interaction between the students' hands and the educational tools is documented, reflecting direct insights into the students' immediate actions. These consequential behaviors are then aligned with the parameters set by the ICAP framework, serving as a robust metric for evaluating student engagement throughout STEM-centric tasks.

This study engages with the SHAP (SHapley Additive exPlanations) methodology, an advanced technique within the realm of interpretable artificial intelligence, to critically analyze the contributory features inherent in the academic performance prediction model. Concurrently, an intriguing observation emerges from the C1 cohort, exhibiting a marginal enhancement in predictive accuracy relative to the established baseline, which is preliminarily set at 50%. This nuanced increment, albeit minimal, signals a critical inference: the interactive dynamics encapsulated within the online classroom environment exert a relatively insubstantial influence on the academic trajectories associated with non-STEM coursework. This revelation underscores the necessity for a differential pedagogical approach, potentially customized to the distinct educational exigencies of STEM and non-STEM curricula.

## IV.    EXPERIMENTAL RESULTS

### A. Evaluation of Emotional Expression of Students

We implemented a quasi-experimental approach to investigate the impact of using RASEDS for recommending adaptive learning materials in STEM education, particularly in enhancing student engagement and self-confidence. This experiment was integrated into the 'Networks Embedded System and Application' course, spanning two academic terms. While students undertook the course independently, collaboration and dialogue were encouraged during the project development phase. Emphasizing IoT and AI, the course required students to leverage their understanding of both software and hardware to devise solutions for real-world challenges, thereby resonating with the fundamental tenets of STEM education (see Fig. 2).

The experimental phase of the study was meticulously structured and spanned duration of seven weeks. This phase was critically segmented into two distinct assessment periods, wherein participants from diverse groups were engaged in comprehensive evaluations. The primary objective of these assessments was to ascertain and quantify two fundamental dimensions: the degree of participant involvement and the perception of personal competence.

At the outset of the experimental phase, in the first week, an initial assessment was conducted. This preliminary evaluation served as a baseline measurement, establishing the initial state of participant engagement and their self-assessed competence. This was crucial for providing a reference point against which any subsequent changes could be measured. The initial assessment was designed to be comprehensive, ensuring that all relevant aspects of involvement and personal competence were adequately captured.

As the program progressed, participants continued their engagement in the designed activities and interventions. This progression was systematically documented and is visually represented in Fig. 3 of the study. The figure illustrates the temporal flow of the program, marking key milestones and the transition from the initial to the final stages.

Fig. 2. STEM education process in Application



Fig. 3. Pretest and posttest engagements of self-efficacy.

In the concluding phase of the program, during the fifth week, a final assessment was conducted. This assessment mirrored the initial one in structure but aimed to capture the evolved state of participant involvement and competence. The comparison between the initial and final assessments was pivotal in determining the effectiveness of the program. It enabled the researchers to quantify the changes in the levels of involvement and personal competence, attributing these changes to the interventions and activities experienced by the participants.

In summary, the experimental phase, with its well-defined and strategically placed assessments, provided a robust framework for evaluating the impact of the program on participant involvement and personal competence. The assessments, anchored at the beginning and end of the program, offered critical insights into the developmental trajectory of the participants, as detailed in Fig. 3.

The outcomes derived from the confusion matrix, presented in Fig. 4, enable the computation of precision, recall, and F1 score for each category of engagement as detected by RASEDS. In the realm of machine learning model evaluation, a confusion matrix serves as a pivotal tool, offering nuanced insights into the classification prowess of a model across various categories [38]. The presented matrix delineates the performance of a classifier in segregating data into five distinct classes: Interactive, Constructive, Active, Passive, and Other. The matrix's structure, with rows representing actual classes and columns depicting predicted classes, provides a comprehensive view of both the model's accuracy and its errors in classification.



Fig. 4. Experimental results.

A closer inspection of the matrix reveals intricate details about the model's performance. For the 'Interactive' class, there is a prominent diagonal element of 81, indicative of a high rate of correct predictions. However, there is a noticeable misclassification with the 'Passive' and 'Constructive' categories, as evidenced by the presence of 8 and 4 instances,

respectively, in these columns. This pattern suggests a certain level of ambiguity or overlap in the defining characteristics of these classes as interpreted by the model.

The 'Constructive' class exhibits impressive prediction accuracy, with 89 instances correctly classified. Nonetheless, there are marginal confusions with the 'Interactive', 'Active', and 'Passive' classes, albeit to a lesser extent than observed in the 'Interactive' class. This points to a generally robust model performances in this category, with room for improvement in differentiating finer nuances between certain classes.

For the 'Active' and 'Passive' classes, the model demonstrates commendable predictive accuracy, as indicated by 91 and 87 correct predictions, respectively. Misclassifications in these categories are relatively lower, suggesting that the model effectively captures the distinct features of these classes. The 'Other' category, with a high correct prediction count of 91, confirms the model's capacity to accurately identify instances that do not conform to the primary classes.

In sum, the confusion matrix provides an invaluable quantitative assessment of the model's classification abilities, highlighting areas of strength and pinpointing aspects that warrant further refinement [39]. Through this detailed analysis, researchers can gain a profound understanding of the model's behavior across varied classifications, guiding targeted improvements in its predictive accuracy.

Intricately woven into this analysis are six pivotal variables, each derived from a comprehensive aggregation of the absolute values of corresponding interactive or emotional metrics within a specific interactive phase. For instance, the variable 'summary_interaction' is computed by summing the absolute SHAP values of various interactive categories during the summary stage, represented formulaically as: summary_interaction = $|ics| + |ims| + |ios| + |ccs| + |cms| + |cos|$. Analogously, 'summary_emotion' encapsulates the emotional undertones of the summary phase, calculated as: summary_emotion = $|ips| + |cps| + |ins| + |cns|$.

## V. DISCUSSION

In this study, a comprehensive literature review was conducted, focusing on the application of machine learning and computer vision techniques across various domains, as cited in references [40-42]. These techniques have been noted for their diverse utility, ranging from healthcare to educational applications. Building on this foundation, the current study specifically applies machine learning and computer vision methods within the realm of STEM education, aiming to explore and expand the educational potential of these innovative technologies.

In the context of STEM education, the integration of computer vision and machine learning (ML) offers transformative potential [43-45]. This research paper has explored how these technologies can be leveraged to enhance self-study methodologies in STEM subjects, with a focus on personalized learning, engagement, and improved learning outcomes.

Personalization of Learning: One of the most significant contributions of ML in STEM education is the ability to tailor educational content to individual students' needs. By analyzing student performance and learning behaviors, ML algorithms can adaptively modify the curriculum, presenting topics in a manner that aligns with each student's unique learning style and pace. This personalization is crucial in self-study environments, where learners often lack the direct guidance of an instructor.

Enhanced Engagement through Computer Vision: The application of computer vision in educational tools has been shown to increase student engagement. By incorporating interactive visual elements and real-time feedback systems, computer vision can make abstract STEM concepts more tangible and comprehensible. This visual interactivity is particularly effective in self-study scenarios, keeping students motivated and engaged in the absence of traditional classroom dynamics.

Data-Driven Insights: ML algorithms provide valuable insights into student learning patterns, identifying areas of difficulty and success. This data can inform the design of future educational content, ensuring that it addresses common challenges and reinforces key concepts. In self-study, these insights become crucial for students to monitor their progress and for educators to understand the efficacy of the learning material.

Overcoming Challenges: Despite the advantages, the integration of computer vision and ML in STEM education is not without challenges. Concerns regarding data privacy, the digital divide, and the need for robust and unbiased algorithms are paramount. Ensuring that these technologies are accessible and equitable for all students, regardless of background or resources, is essential for their successful implementation in educational settings.

Future Directions: Looking forward, the continued development and refinement of ML and computer vision technologies promise even greater advancements in STEM education. The potential integration of augmented reality (AR) and virtual reality (VR) technologies, combined with ML-driven personalized learning paths, could revolutionize the way STEM subjects are taught and learned. Additionally, the ongoing improvement of algorithmic transparency and fairness will be crucial in ensuring that these technologies serve all students effectively.

In conclusion, the application of computer vision and machine learning in STEM-education self-study represents a significant step forward in educational technology. These tools offer the potential for highly personalized, engaging, and effective learning experiences. However, careful attention must be paid to the challenges and ethical considerations that come with the implementation of such advanced technologies. As the field progresses, continuous evaluation and adaptation will be necessary to fully realize the benefits of these innovations in STEM education.

## VI. Conclusion

This study's journey into the realms of advanced technology's application within educational settings, particularly through the Real-time Automated STEM Engagement Detection System (RASEDS), unveils new horizons in STEM-related pedagogies. The evidence presented underscores the potential of such innovative intersections between technology and education, where systems like RASEDS are not mere analytical tools but catalysts for transformative educational experiences. The ability of RASEDS to discern engagement levels accurately heralds a future where learning can be genuinely individualized, responding in real-time to students' engagement fluctuations. Moreover, the observed enhancement in self-efficacy among learners signals a profound impact on learners' psychological resources, potentially influencing their academic trajectories and career paths in STEM fields.

However, the journey does not conclude here. While the findings affirm the positive trajectories, they also cast light on the complexities and multi-dimensional challenges within technology-integrated education. Future research needs to navigate these sophisticated dynamics, including the nuanced understanding of engagement and self-efficacy, ethical considerations surrounding data security, and the psychological safety of learners. Furthermore, pedagogical strategies must evolve in tandem with these technological advancements, ensuring that human-centric learning remains at the core of educational endeavors. As we stand on the brink of this new era, the responsibility is collective—educators, technologists, policymakers, and researchers must collaborate to ensure these innovations are harnessed responsibly, ethically, and with the holistic development of learners in mind.

### References

[1] Alabdulhadi, A., & Faisal, M. (2021). Systematic literature review of STEM self-study related ITSs. Education and Information Technologies, 26(2), 1549-1588.

[2] Xu, W., & Ouyang, F. (2022). The application of AI technologies in STEM education: a systematic review from 2011 to 2021. International Journal of STEM Education, 9(1), 1-20.

[3] Tang, J., Zhou, X., Wan, X., & Ouyang, F. (2022). A Systematic Review of AI Applications in Computer-Supported Collaborative Learning in STEM Education. Artificial Intelligence in STEM Education: The Paradigmatic Shifts in Research, Education, and Technology, 333.

[4] Lämsä, J., Virtanen, A., Tynjälä, P., Maunuksela, J., & Koskinen, P. (2023). Exploring students' perceptions of self-assessment in the context of problem solving in STEM. LUMAT, 11(2).

[5] Yurchenko, A., Yurchenko, K., Proshkin, V., & Semenikhina, O. (2022). World Practices of STEM Education Implementation: Current Problems and Results. International Journal of Research in E-learning, 8(2), 1-20.

[6] B. Omarov, S. Narynov, Z. Zhumanov, A. Gumar and M. Khassanova, "A skeleton-based approach for campus violence detection," Computers, Materials & Continua, vol. 72, no.1, pp. 315–331, 2022.

[7] Omarov, B., Omarov, B., Shekerbekova, S., Gusmanova, F., Oshanova, N., Sarbasova, A., ... & Sultan, D. (2019). Applying face recognition in video surveillance security systems. In Software Technology: Methods and Tools: 51st International Conference, TOOLS 2019, Innopolis, Russia, October 15–17, 2019, Proceedings 51 (pp. 271-280). Springer International Publishing.

[8] Mystakidis, S., Papantzikos, G., & Stylios, C. (2021, September). Virtual reality escape rooms for STEM education in industry 4.0: Greek teachers perspectives. In 2021 6th South-East Europe Design Automation, Computer Engineering, Computer Networks and Social Media Conference (SEEDA-CECNSM) (pp. 1-5). IEEE.

[9] Zhang, X., Zhang, B., & Zhang, F. (2023). Student-centered case-based teaching and online–offline case discussion in postgraduate courses of

computer science. International Journal of Educational Technology in Higher Education, 20(1), 6.

[10] Omarov, B., Suliman, A., & Kushibar, K. (2016). Face recognition using artificial neural networks in parallel architecture. Journal of Theoretical and Applied Information Technology, 91(2), 238.

[11] Tuong, D. H., Tran, T. B., & Nguyen, D. D. (2023). Digital Transformation for Vietnam Education: From Policy to School Practices. In New Challenges and Opportunities in Physics Education (pp. 293-311). Cham: Springer Nature Switzerland.

[12] Gull, R. A., Izham, M. D. B. M., & Qadir, J. (2023, May). 1 Robotics Primer for Independent Learners: Background, Curriculum, Resources, and Tips. In 2023 IEEE Global Engineering Education Conference (EDUCON) (pp. 1-9). IEEE.

[13] Hlukhaniuk, V., Solovei, V., Tsvilyk, S., & Shymkova, I. (2020, May). STEAM education as a benchmark for innovative training of future teachers of labour training and technology. In SOCIETY. INTEGRATION. EDUCATION. Proceedings of the International Scientific Conference (Vol. 1, pp. 211-221).

[14] Christensen, D., Singelmann, L., Sleezer, R., & Siverling, E. A. (2023, June). A Self-Study of Faculty Methods, Attitudes, and Perceptions of Oral Engineering Exams. In 2023 ASEE Annual Conference & Exposition.

[15] Tong, D. H., Uyen, B. P., & Ngan, L. K. (2022). The effectiveness of blended learning on students' academic achievement, self-study skills and learning attitudes: A quasi-experiment study in teaching the conventions for coordinates in the plane. Heliyon, 8(12).

[16] Horng-Jyh, P. W., Cheng, C. H. K., Tah, B. L. Y., Lie, T. H., Beng, J. S. T., Guan, R. O. B., ... & Yongqing, Z. (2023, June). Virtual Lab Workspace for Programming Computers–Towards Agile STEM Education. In International Workshop on Learning Technology for Education Challenges (pp. 55-68). Cham: Springer Nature Switzerland.

[17] Narynov, S., Zhumanov, Z., Gumar, A., Khassanova, M., & Omarov, B. (2021, October). Chatbots and Conversational Agents in Mental Health: A Literature Review. In 2021 21st International Conference on Control, Automation and Systems (ICCAS) (pp. 353-358). IEEE.

[18] Yao, X. (2022). Design and research of artificial intelligence in multimedia intelligent question-answering system and self-test system. Advances in Multimedia, 2022.

[19] Lim, J., Jo, H., Zhang, B. T., & Park, J. (2021). Passive versus active: Frameworks of active learning for linking humans to machines. In Proceedings of the Annual Meeting of the Cognitive Science Society (Vol. 43, No. 43).

[20] Xue, H. (2022). A new integrated teaching mode for labor education course based on STEAM education. International Journal of Emerging Technologies in Learning (iJET), 17(2), 128-142.

[21] Fung, C. H., Poon, K. K., & Ng, S. P. (2022). Fostering student teachers' 21st century skills by using flipped learning by teaching in STEM education. Eurasia Journal of Mathematics, Science and Technology Education, 18(12), em2204.

[22] Tinh, P. T., Duc, N. M., Yuenyong, C., Kieu, N. T., & Nguyen, T. T. (2021, March). Development of STEM education learning unit in context of Vietnam Tan Cuong Tea village. In Journal of Physics: Conference Series (Vol. 1835, No. 1, p. 012060). IOP Publishing.

[23] Cooper, G. (2023). Examining science education in chatgpt: An exploratory study of generative artificial intelligence. Journal of Science Education and Technology, 32(3), 444-452.

[24] Horng-Jyh, P. W., Cheng, C. H. K., Tah, B. L. Y., Lie, T. H., Beng, J. S. T., Guan, R. O. B., ... & Yongqing, Z. (2023, June). Towards Agile STEM Education. In Learning Technology for Education Challenges: 11th International Workshop, LTEC 2023, Bangkok, Thailand, July 24–27, 2023, Proceedings (p. 55). Springer Nature.

[25] Nguyen, H. D., Tran, T. V., Pham, X. T., Huynh, A. T., Pham, V. T., & Nguyen, D. (2022). Design intelligent educational chatbot for information retrieval based on integrated knowledge bases. IAENG International Journal of Computer Science, 49(2), 531-541.

[26] Wu, T. T., Lin, C. J., Pedaste, M., & Huang, Y. M. (2023, August). The Effect of Chatbot Use on Students' Expectations and Achievement in STEM Flipped Learning Activities: A Pilot Study. In International

[27] Sultanovich, O. B., Ergeshovich, S. E., Duisenbekovich, O. E., Balabekovna, K. B., Nagashbek, K. Z., & Nurlakovich, K. A. (2016). National Sports in the Sphere of Physical Culture as a Means of Forming Professional Competence of Future Coach Instructors. Indian Journal of Science and Technology, 9(5), 87605-87605.

[28] Qian, Y., Hambrusch, S., Yadav, A., & Gretter, S. (2018). Who needs what: Recommendations for designing effective online professional development for computer science teachers. Journal of Research on Technology in Education, 50(2), 164-181.

[29] Yadav, A., Connolly, C., Berges, M., Chytas, C., Franklin, C., Hijón-Neira, R., ... & Warner, J. R. (2022). A review of international models of computer science teacher education. Proceedings of the 2022 Working Group Reports on Innovation and Technology in Computer Science Education, 65-93.

[30] Lin, X., Liu, H., Sun, Q., Li, X., Qian, H., Sun, Z., & Lam, T. L. (2022). Applying project-based learning in artificial intelligence and marine discipline: An evaluation study on a robotic sailboat platform. IET Cyber-Systems and Robotics, 4(2), 86-96.

[31] UmaMaheswaran, S. K., Prasad, G., Omarov, B., Abdul-Zahra, D. S., Vashistha, P., Pant, B., & Kaliyaperumal, K. (2022). Major challenges and future approaches in the employment of blockchain and machine learning techniques in the health and medicine. Security and Communication Networks, 2022.

[32] Saraswathi, K., Devadharshini, B., Kavina, S., & Srinidhi, S. (2023, February). Prediction on Impact of Electronic Gadgets in Students Life using Machine Learning. In 2023 7th International Conference on Computing Methodologies and Communication (ICCMC) (pp. 340-345). IEEE.

[33] Vdovinskienė, S. (2023). Using Flipped Classroom as an Active Teaching Method for Teaching Engineering Graphics. Baltic Journal of Modern Computing, 11(3).

[34] Zufarova, O., Kondratieva, V., & Zhirosh, O. (2021, November). Learning Environment-What Matters for the High Ability Computer Science Students?. In 2021 World Engineering Education Forum/Global Engineering Deans Council (WEEF/GEDC) (pp. 144-152). IEEE.

[35] Theissler, A., & Ritzer, P. (2022, March). EduML: An explorative approach for students and lecturers in machine learning courses. In 2022 IEEE Global Engineering Education Conference (EDUCON) (pp. 921-928). IEEE.

[36] Ormanci, Ü. (2020). Thematic content analysis of doctoral theses in STEM education: Turkey context. Journal of Turkish Science Education, 17(1), 126-146.

[37] Scull, W. R., Perkins, M. A., Carrier, J. W., & Barber, M. (2023). Community college institutional researchers' knowledge, experience, and perceptions of machine learning. Community College Journal of Research and Practice, 47(5), 354-368.

[38] Xiaohong, C., Jie, L., Zhibin, M., & Li, X. (2021, June). Teaching research on the cultivation of computational thinking ability by using information technology. In 2021 2nd International Conference on Artificial Intelligence and Education (ICAIE) (pp. 564-567). IEEE.

[39] Corlu1, M. S., Kurutas, B. S., & Ozel, S. (2023). Effective online professional development: A facilitator's perspective. In Research On STEM Education in the Digital Age: Proceedings of the ROSEDA Conference (Vol. 6, p. 9). WTM-Verlag Münster.

[40] Rusnilawati, R., Ali, S. R. B., Hanapi, M. H. M., Sutama, S., & Rahman, F. (2023). The Implementation of Flipped Learning Model and STEM Approach in Elementary Education: A Systematic Literature Review. European Journal of Educational Research, 12(4).

[41] Omarov, B., Altayeva, A., Suleimenov, Z., Im Cho, Y., & Omarov, B. (2017, April). Design of fuzzy logic based controller for energy efficient operation in smart buildings. In 2017 First IEEE International Conference on Robotic Computing (IRC) (pp. 346-351). IEEE.

[42] Feng, H., & Wang, J. (2022). Learning in a Digital World: Perspective on Interactive Technologies for Formal and Informal Education: edited by Paloma Díaz, Andri Ioannou, Kaushal Kumar Bhagat, and J. Michael Spector, Springer (Singapore, Singapore), 2019, xviii+ 339pp.,€ 96.29 (ebook), ISBN 978-981-13-8265-9.

[43] Aboamer, M. A., Sikkandar, M. Y., Gupta, S., Vives, L., Joshi, K., Omarov, B., & Singh, S. K. (2022). An investigation in analyzing the food quality well-being for lung cancer using blockchain through cnn. Journal of Food Quality, 2022.

[44] Katal, A., Upadhyay, J., & Singh, V. K. (2023). Blended learning in COVID-19 era and way-forward. In Sustainable Blended Learning in STEM Education for Students with Additional Needs (pp. 55-85). Singapore: Springer Nature Singapore.

[45] Kazmi, H., Munné-Collado, Í., Tidriri, K., Nordström, L., Gielen, F., & Driesen, J. (2022). Data science and energy: some lessons from Europe on higher education course design and delivery. Harvard Data Science Review, 4(1).

# Automated Fruit Sorting in Smart Agriculture System: Analysis of Deep Learning-based Algorithms

Cheng Liu[1], Shengxiao Niu[2]*

College of Digital Business, Jiangsu Vocational Institute of Commerce, Nanjing 210000, Jiangsu, China[1]
Handan Polytechnic College, Handan 056000, Hebei, China[2]

*Abstract*—**Automated fruit sorting plays a crucial role in smart agriculture, enabling efficient and accurate classification of fruits based on various quality parameters. Traditionally, rule-based and machine-learning methods have been employed for fruit sorting, but in recent years, deep learning-based approaches have gained significant attention. This paper investigates deep learning methods for fruit sorting and justifies their prevalence in the field. Therefore, it is necessary to address these limitations and improve the effectiveness of CNN-based fruit sorting methods. This research paper presents a comprehensive analysis of CNN-based methods, highlighting their strengths and limitations. This analysis aims to contribute to advancing automated fruit sorting in smart agriculture and provide insights for future research and development in deep learning-based fruit sorting techniques.**

*Keywords*—*Smart agriculture; automated fruit sorting; deep learning; Convolutional Neural Network (CNN); analysis*

## I. INTRODUCTION

Automated fruit sorting has emerged as a promising technology in the field of smart agriculture, revolutionizing the way fruits are cultivated, harvested, and processed [1, 2]. These technologies integrate advanced sensing, data analytics, and automation techniques to improve productivity, efficiency, and quality in fruit production and processing [3, 4]. Automated fruit sorting plays a vital role in the post-harvest stage, ensuring accurate and efficient classification of fruits based on various quality parameters such as size, color, shape, and ripeness [5]. One of the key components of automated fruit sorting systems is video-based fruit sorting [6], which utilizes computer vision and image processing techniques to analyze visual information and make real-time applications [7, 8].

Video-based fruit sorting has gained significant attention due to its non-destructive nature, high-speed operation, and ability to handle a large volume of fruits [9, 10]. This approach involves capturing video footage of fruits from multiple angles and utilizing computer vision algorithms to extract relevant features for classification. By analyzing the visual characteristics of fruits, video-based sorting systems can accurately classify them into different categories, ensuring consistent quality and reducing manual labor.

In recent years, there have been remarkable advancements in video-based fruit sorting systems driven by the rapid progress in computer vision, machine learning, and deep learning techniques [10, 11]. These technologies have enabled the development of more sophisticated and efficient algorithms for fruit classification, leading to improved accuracy and speed

in sorting operations [12, 13]. Deep learning, in particular, has shown great potential in fruit sorting applications, leveraging its ability to learn discriminative features from large-scale data automatically.

Deep learning-based approaches have demonstrated superior performance in various computer vision tasks, and fruit sorting is no exception [14, 15]. Convolutional Neural Networks (CNNs) have emerged as a popular choice for fruit classification due to their ability to extract hierarchical features from images [10, 16, 17]. The context of the CNNs, significant features are automatically learned during the training process [18]. CNNs automatically identify and extract relevant patterns and features from input data through convolutional layers, optimizing the network's parameters to minimize the difference between predicted and actual target labels. Additionally, Recurrent Neural Networks (RNNs) and hybrid models combining CNNs and RNNs have been explored to capture spatial and temporal information in video-based fruit sorting.

Although significant progress has been made in deep learning-based fruit sorting, there are still some limitations and research gaps that need to be addressed. Firstly, the lack of annotated datasets specific to fruit sorting poses challenges for training and evaluating deep learning models. Secondly, the generalization and robustness of deep learning models across different fruit types and environmental conditions need to be investigated further. Finally, the computational complexity and deployment feasibility of deep learning models in real-world fruit sorting systems requires careful consideration.

Therefore, this review paper aims to address these research gaps and present an in-depth investigation and analysis of deep learning methods for fruit sorting. By conducting this investigation, we aim to shed light on the potential of deep learning methods in improving the efficiency, accuracy, and scalability of automated fruit sorting systems, contributing to the advancement of smart agriculture and post-harvest technologies. The contributions of this study are three-fold:

*1) A* comprehensive review of the most recent deep learning-based approaches for fruit sorting, highlighting their strengths and limitations;

*2) An* analysis of current research gaps and challenges in CNN-based methods;

*3) Addressing* potential strategies and future directions to overcome these challenges and advance the field of deep learning-based fruit sorting.

## II. RELATED WORKS

This section provides related works focusing on grading and sorting fruit using machine learning and deep learning-based approaches.

Patil et al. [19] focus on the grading and sorting technique of dragon fruits using machine learning algorithms. The study explores the application of machine learning algorithms, specifically support vector machine (SVM) and random forest, for grading and sorting dragon fruits based on their quality attributes. The findings demonstrate that both SVM and random forest models achieved high accuracy in classifying the dragon fruits into different grades. However, the study also highlights certain limitations, such as the need for a large and diverse dataset to improve the models' performance and the challenges of integrating the grading and sorting system into an automated production line. This research contributes to the development of efficient grading and sorting techniques for dragon fruits using machine-learning algorithms while also acknowledging the areas that require further exploration and improvement.

Gill and Khehra [20] focused on fruit image classification using deep learning techniques. The study aims to develop an accurate and efficient system for automatically classifying fruits based on their images. The researchers employ deep learning models, such as convolutional neural networks (CNNs), to extract meaningful features from fruit images and train classification models. The findings demonstrate the effectiveness of deep learning in accurately identifying different types of fruits, achieving high classification accuracy. The proposed system has practical applications in fruit sorting and quality control processes, enabling faster and more reliable classification compared to traditional methods. Overall, this research contributes to the field of automated fruit classification using deep learning, showcasing the potential of this approach in various fruit-related industries.

Kumar and Parkavi [21] provided a comprehensive review of the quality grading of fruits and vegetables using image processing techniques and machine learning. The study examines various image processing methods, such as color analysis, texture analysis, and shape analysis, and discusses their applications in assessing the quality attributes of fruits and vegetables. Machine learning algorithms, including support vector machine (SVM), random forest, and artificial neural networks (ANN), are investigated for automated quality grading. The findings highlight the effectiveness of image processing techniques coupled with machine learning in accurately grading fruits and vegetables based on their quality parameters. However, the paper also recognizes certain limitations, such as the need for robust and diverse datasets, standardized grading criteria, and real-time implementation challenges. This review serves as a valuable resource for researchers and practitioners in the field of automated fruit and vegetable quality grading while emphasizing the areas that require further research and development to overcome the existing limitations.

Chakraborty et al. [22] presented the development of a real-time automatic citrus fruit grading and sorting machine using a computer vision-based adaptive deep learning model. The study aims to improve the efficiency and accuracy of citrus fruit grading by leveraging advanced machine-learning techniques. The findings demonstrate that the proposed system, equipped with an optimized deep learning model, achieves high accuracy in grading citrus fruits based on quality attributes such as size, color, and shape. The system effectively handles various challenges encountered in citrus fruit grading, such as variations in fruit appearance and lighting conditions. However, the paper also acknowledges certain limitations, including the need for a large and diverse dataset to enhance the model's performance further. This research contributes to the development of a practical and efficient citrus fruit grading system while highlighting the potential for further advancements and improvements in deep learning-based approaches for fruit sorting applications.

## III. METHODOLOGY

With the continuous advancements in Convolutional Neural Network (CNN) architectures and the availability of well-annotated fruit datasets, CNN-based frameworks have emerged as valuable tools for automating fruit sorting processes across various industries, including agriculture, food processing, and packaging.

In this research study, we focus on the evaluation and analysis of existing CNN-based approaches for fruit disease detection. We specifically investigate the performance of popular CNN frameworks, namely DenseNet, InceptionV3, ResNet, VGGNet, Xception, MobileNet, NASNet, EfficientNet, and SqueezeNet. To achieve this, we conduct extensive experiments using these models and collect the resulting performance metrics. In addition to our experiments, we gather data from previously published research works to augment our analysis. We extract performance measurements such as sensitivity, specificity, and accuracy from these studies. By incorporating a diverse range of sources, we aim to provide a comprehensive overview of the effectiveness of CNN-based approaches in fruit disease detection. For the dataset, this study uses Fruits 360. The Fruits 360 is a large-scale dataset of images containing fruits and vegetables, which can be used for various computer vision tasks such as classification, segmentation, and detection. The dataset consists of 90380 images of 131 different types of fruits and vegetables, with each image having a size of 100x100 pixels.

### A. CNN based Methods

This study focuses on exploring and analyzing the effectiveness of CNN-based approaches for automated fruit sorting. Extensive experiments are conducted to evaluate the performance of various models, and the results are carefully gathered and analyzed. Additionally, valuable insights are gathered from previously published research works, where performance measurements based on sensitivity, specificity, and accuracy metrics are collected and compared. By examining experimental findings and existing literature, this study aims to provide comprehensive insights into the effectiveness and potential of CNN-based methods for automated fruit-sorting applications.

*1) ResNet:* ResNet, short for Residual Neural Network, is a deep learning architecture that revolutionized image

classification tasks, including automated fruit sorting [23]. ResNet introduces skip connections that allow the network to learn residual mappings, making it easier to train deeper networks [24]. This architecture helps in overcoming the degradation problem in very deep networks and enables the accurate classification of fruits based on their visual characteristics. Fig. 1 shows the structure of the ResNet model.

*2) InceptionV3:* InceptionV3 is a widely used deep convolutional neural network architecture for automated fruit sorting. It employs a combination of 1x1, 3x3, and 5x5 convolutional filters to capture various scales of features in the input images [26]. InceptionV3's inception modules efficiently capture both local and global patterns, allowing for accurate classification and identification of fruit types (see Fig. 2).

*3) VGGNet:* VGGNet is a classic deep convolutional neural network architecture that has been applied to automated fruit sorting. It consists of multiple convolutional layers with small receptive fields, followed by fully connected layers [28]. VGGNet's uniform architecture and deeper network depth allow it to capture intricate visual features, leading to robust fruit classification and sorting capabilities (see Fig. 3).

*4) DenseNet:* DenseNet is another deep learning architecture commonly utilized in fruit sorting tasks. DenseNet introduces dense connections, where each layer is directly connected to every other layer in a feed-forward fashion [30]. These dense connections enable feature reuse and encourage gradient flow, resulting in more efficient and accurate classification of fruits based on their attributes (see Fig. 4).

*5) MobileNet:* MobileNet is a lightweight deep-learning architecture designed for mobile and resource-constrained devices. It employs depth-wise separable convolutions to reduce the computational cost while preserving accuracy [32]. MobileNet-based models are efficient for fruit sorting applications where computational resources are limited (see Fig. 5).

*6) NASNet:* NASNet, short for Neural Architecture Search Network, is an architecture discovered using neural architecture search techniques. It automatically searches for optimal network architectures for fruit sorting, resulting in highly efficient and accurate models [34]. NASNet-based models can adapt to different fruit sorting tasks by automatically learning the optimal network structure (see Fig. 6).



Fig. 1. The structure of ResNet [25].



Fig. 2. Inception V3 structure [27].

Fig. 3. The structure of VGGNet [29].



Fig. 4. The structure of DenseNet [31].



Fig. 5. The structure of MobileNet [33].

Fig. 6.  Two structures of NASNet [35].



Fig. 7.  The structure of EfficientNet [36].



Fig. 8.  The structure of SqueezeNet [39].

*7) EfficientNet:* EfficientNet is a family of deep learning models that achieve state-of-the-art performance with significantly fewer parameters and computational resources. These models employ a compound scaling method that balances model depth, width, and resolution to achieve optimal performance [15, 34]. EfficientNet-based models provide excellent accuracy and efficiency for automated fruit sorting tasks (see Fig. 7).

*8) SqueezeNet:* SqueezeNet is a lightweight deep-learning architecture that achieves high accuracy with a reduced number of parameters. It utilizes fire modules, which consist of both 1x1 and 3x3 filters, to efficiently capture and process fruit image features [37, 38]. SqueezeNet is particularly suitable for resource-constrained environments while maintaining competitive performance in fruit sorting (see Fig. 8).

### B. Performance Measurements

In CNN-based fruit sorting, performance measurements such as sensitivity, specificity, and accuracy play a crucial role in evaluating the effectiveness of the models. These metrics are derived from the concepts of True Positive (TP), False Negative (FN), True Negative (TN), and False Positive (FP). The definitions of these metrics are as follows:

True Positive (TP): It represents the number of correctly classified positive instances, i.e., the number of diseased fruits correctly identified by the CNN model.

False Negative (FN): It refers to the number of positive instances that were incorrectly classified as negative, i.e., the number of diseased fruits that were wrongly identified as healthy or undetected by the CNN model.

True Negative (TN): It represents the number of correctly classified negative instances, i.e., the number of healthy fruits correctly identified by the CNN model.

False Positive (FP): It refers to the number of negative instances that were incorrectly classified as positive, i.e., the number of healthy fruits that were wrongly identified as diseased by the CNN model.

Based on these definitions, we can calculate the following performance measurements, Sensitivity (True Positive Rate or Recall):

Sensitivity measures the proportion of correctly classified positive instances out of all the actual positive instances. It indicates the model's ability to detect and classify diseased fruits accurately.

Specificity (True Negative Rate): Specificity measures the proportion of correctly classified negative instances out of all the actual negative instances. It evaluates the model's ability to accurately identify healthy fruits without misclassifying them as diseased.

Accuracy: Accuracy represents the overall correctness of the model's predictions by calculating the proportion of correctly classified cases, positive and negative, out of the total number of cases. The corresponding equations for sensitivity, specificity and accuracy are as follows:

$$Sensitivity = TP / (TP + FN)$$

$$Specificity = TN / (TN + FP)$$

$$Accuracy = (TP + TN) / (TP + TN + FP + FN)$$

### IV. ANALYSIS OF CNN-BASED METHODS

In this section, a performance analysis of CNN-based frameworks is presented. We have chosen the widely adopted CNN frameworks, namely DenseNet, InceptionV3, ResNet, VGGNet, Xception, MobileNet, NASNet, EfficientNet, and SqueezeNet models. To thoroughly investigate their performance, we conducted a series of comprehensive experiments and meticulously collected the corresponding results. Additionally, we gathered relevant data from previously published research works, which provided valuable insights into the models' performance. The performance measurements were evaluated using sensitivity, specificity, and accuracy metrics, enabling a robust assessment of the models' capabilities. Specificity, sensitivity, accuracy, and associated metrics such as true positive (TP), true negative (TN), false positive (FP), and false negative (FN) are considered as the most popular and fundamental metrics for technical analysis and performance measurement, particularly in classification tasks such as automated fruit sorting using CNN-based models.

The sensitivity values provide insights into how well each CNN-based method can identify diseased fruits, a critical aspect of fruit sorting for disease detection. Sensitivity is particularly relevant in applications where minimizing false negatives is essential, ensuring that diseased fruits are not overlooked.

As shown in Fig. 9, we observe that EfficientNet demonstrates the highest sensitivity value of 0.93, indicating its strong capability to identify diseased fruits accurately. InceptionV3 follows closely with a sensitivity of 0.92, highlighting its effectiveness in detecting diseased instances. Xception and ResNet also show notable sensitivity values of 0.91 and 0.9, respectively.

DenseNet and NASNet have sensitivity values of 0.89 and 0.9, respectively, indicating their ability to capture most of the diseased fruits but with a slightly lower performance compared to the aforementioned methods. VGGNet has a sensitivity of 0.88, while MobileNet and SqueezeNet have lower sensitivities of 0.87 and 0.85, respectively.

Based on these sensitivity values, it can be inferred that EfficientNet, InceptionV3, Xception, and ResNet exhibit relatively higher performance in correctly identifying diseased fruits. These models are likely to be more reliable in fruit disease detection applications.

Fig. 9.   Analysis of CNN methods based on sensitivity metric.



Fig. 10.  Analysis of CNN methods based on specificity metric.

In fruit sorting applications, a high specificity value is desirable as it ensures that healthy fruits are correctly recognized and avoids misclassifying them as diseased. These specificity values help in evaluating the performance of different CNN-based methods and can guide the selection of appropriate models for fruit sorting tasks. These specificity values represent the proportion of correctly classified negative cases (healthy fruits) out of all the actual negative cases. A higher specificity value indicates a better ability of the model to accurately identify healthy fruits without misclassifying them as diseased.

Based on the specificity values provided in Fig. 10, we can observe that EfficientNet has the highest specificity (0.94),

followed by InceptionV3, ResNet, Xception, and NASNet, which all have a specificity of 0.92. DenseNet, VGGNet, and MobileNet have a specificity of 0.91, while SqueezeNet has the lowest specificity at 0.88. These specificity values provide insights into the models' performance in accurately identifying healthy fruits in the fruit sorting process. A higher specificity indicates a lower chance of misclassifying healthy fruits as diseased, which is desirable for efficient fruit sorting applications. Therefore, among the methods listed, EfficientNet stands out with the highest specificity value of 0.94, indicating its strong capability to identify healthy fruits accurately. InceptionV3, ResNet, Xception, and NASNet also exhibit high specificity values of 0.92, highlighting their effectiveness in correctly classifying healthy fruits.

Fig. 11. Analysis of CNN methods based on accuracy rate.

Inaccuracy measurement, Fig. 11 presents accuracy rates for different CNN-based methods used in a certain application. Accuracy measures the overall correctness of a model's predictions and represents the proportion of correctly classified instances (both positive and negative) out of the total instances.

By analyzing the accuracy results, we observe that EfficientNet achieves the highest accuracy rate of 0.95, indicating its strong performance in accurately classifying diseased and healthy fruits. InceptionV3 follows closely with an accuracy of 0.94, suggesting its effectiveness in achieving correct predictions. ResNet, and Xception, models also exhibit high accuracy values of 0.93, highlighting their reliability in fruit classification tasks.

DenseNet, VGGNet, and NASNet demonstrate an accuracy of 0.92, indicating similar performance in achieving correct classifications. MobileNet, with an accuracy of 0.91, performs slightly lower than the aforementioned methods. SqueezeNet, however, shows a lower accuracy of 0.89, suggesting it may not perform as well in accurately classifying fruit instances.

Based on these accuracy values, EfficientNet stands out as the top-performing model, closely followed by InceptionV3, ResNet, and Xception. These models have demonstrated a higher capability to achieve accurate predictions and can be considered reliable choices for fruit classification tasks.

As results, DenseNet's strength lies in its effective feature reuse and alleviation of the vanishing gradient problem through dense connectivity, enhancing parameter efficiency. InceptionV3 excels at capturing multi-scale features with its inception modules, suitable for diverse object recognition tasks, but its complex architecture may lead to longer training times. ResNet introduces residual connections, enabling the training of very deep networks, but its increased complexity may demand higher computational resources. VGGNet, with its simple and uniform architecture, performs well on image recognition tasks but is susceptible to overfitting. Xception efficiently employs depth-wise separable convolutions, though

it may require more training data. MobileNet, designed for mobile and edge devices, balances accuracy and efficiency but may lack representation capacity. NASNet's use of neural architecture search enhances performance but demands significant computational resources. EfficientNet achieves high accuracy with improved parameter efficiency but may be computationally expensive to train. SqueezeNet's compact design prioritizes parameter efficiency for edge devices but may sacrifice some accuracy. These nuances in strengths and limitations provide insights into the trade-offs associated with each CNN-based framework, aiding informed choices for fruit sorting applications in smart agriculture.

## V. CONCLUSION

This study emphasizes the importance of automated fruit sorting in smart agriculture and the growing significance of deep learning-based approaches in comparison to traditional methods. The need to address limitations in Convolutional Neural Network (CNN)-based fruit sorting methods is acknowledged, prompting the research paper to conduct a comprehensive analysis. This analysis aims to highlight both the strengths and limitations of CNN-based methods for fruit sorting, with the overarching goal of advancing automated fruit sorting in smart agriculture. By focusing on these features, the paper aims to provide valuable insights for future research and development, contributing to the continual improvement of deep learning-based fruit sorting techniques in the agricultural domain. This research study investigates the use of CNN-based frameworks for automating fruit sorting detection in industries such as agriculture, food processing, and packaging. The study focuses on evaluating popular CNN models, including DenseNet, InceptionV3, ResNet, VGGNet, Xception, MobileNet, NASNet, EfficientNet, and SqueezeNet. Through extensive experiments and analysis of performance metrics such as sensitivity, specificity, and accuracy, the study aims to provide a comprehensive understanding of the strengths and limitations of these models. The findings will contribute to the development of more accurate and reliable systems for fruit

sorting algorithms in agriculture, leading to improved efficiency and productivity in the industry. Directions for future works that can be pursued as investigating the effectiveness of ensemble methods in improving the performance of fruit sorting algorithms. By exploring ensemble techniques such as bagging, boosting, or stacking, researchers can examine how the combination of multiple CNN models can further improve the fruit sorting process. Another future work will focus on the real-time implementation and deployment of the CNN-based fruit sorting algorithms. While the current study evaluates the performance of different CNN models, their practical application in real-time sorting systems is an important aspect that requires further exploration.

## REFERENCES

[1] H. M. T. Abbas, U. Shakoor, M. J. Khan, M. Ahmed, and K. Khurshid, "Automated sorting and grading of agricultural products based on image processing," in 2019 8th international conference on information and communication technologies (ICICT), 2019: IEEE, pp. 78-81.

[2] K. K. Paul et al., "Smart Agriculture Using UAV and Deep Learning: A Systematic Review," Internet of Things, pp. 1-16, 2022.

[3] E. Mavridou, E. Vrochidou, G. A. Papakostas, T. Pachidis, and V. G. Kaburlasos, "Machine vision systems in precision agriculture for crop farming," Journal of Imaging, vol. 5, no. 12, p. 89, 2019.

[4] H. Tian, T. Wang, Y. Liu, X. Qiao, and Y. Li, "Computer vision technology in agricultural automation—A review," Information Processing in Agriculture, vol. 7, no. 1, pp. 1-19, 2020.

[5] N. Janu and A. Kumar, "Automated fruit grading system using image fusion," in Smart agricultural services using deep learning, big data, and IoT: IGI Global, 2021, pp. 32-45.

[6] H. Basri, I. Syarif, S. Sukaridhoto, and M. F. Falah, "Intelligent system for automatic classification of fruit defect using faster region-based convolutional neural network (faster r-CNN)," Jurnal Ilmiah Kursor, vol. 10, no. 1, 2019.

[7] A. A. Mei Choo Ang, Kok Weng Ng, Elankovan Sundararajan, Marzieh Mogharrebi, Teck Loon Lim, "Multi-core Frameworks Investigation on A Real-Time Object Tracking Application," Journal of Theoretical & Applied Information Technology, 2014.

[8] H. Kang and C. Chen, "Fast implementation of real-time fruit detection in apple orchards using deep learning," Computers and Electronics in Agriculture, vol. 168, p. 105108, 2020.

[9] H.-K. Wu, J.-S. Wang, and Y.-H. Chen, "Development of fruit grading system based on image recognition," in 2020 IEEE 2nd international conference on architecture, construction, environment and hydraulics (ICACEH), 2020: IEEE, pp. 26-27.

[10] A. Nasiri, A. Taheri-Garavand, and Y.-D. Zhang, "Image-based deep learning automated sorting of date fruit," Postharvest biology and technology, vol. 153, pp. 133-141, 2019.

[11] M. Sugadev, K. Sucharitha, I. R. Sheeba, and B. Velan, "Computer vision based automated billing system for fruit stores," in 2020 Third International Conference on Smart Systems and Inventive Technology (ICSSIT), 2020: IEEE, pp. 1337-1342.

[12] S. K. Behera, A. K. Rath, A. Mahapatra, and P. K. Sethy, "Identification, classification & grading of fruits using machine learning & computer intelligence: a review," Journal of Ambient Intelligence and Humanized Computing, pp. 1-11, 2020.

[13] A. Bhargava and A. Bansal, "Machine learning based quality evaluation of mono-colored apples," Multimedia Tools and Applications, vol. 79, pp. 22989-23006, 2020.

[14] I. M. Nasir et al., "Deep learning-based classification of fruit diseases: An application for precision agriculture," Comput. Mater. Contin, vol. 66, no. 2, pp. 1949-1962, 2021.

[15] L. T. Duong, P. T. Nguyen, C. Di Sipio, and D. Di Ruscio, "Automated fruit recognition using EfficientNet and MixNet," Computers and Electronics in Agriculture, vol. 171, p. 105326, 2020.

[16] Z. Zhou et al., "Advancement in artificial intelligence for on-farm fruit sorting and transportation," Frontiers in Plant Science, vol. 14, p. 1082860, 2023.

[17] X. Zhang, Y. Xun, and Y. Chen, "Automated identification of citrus diseases in orchards using deep learning," Biosystems Engineering, vol. 223, pp. 249-258, 2022.

[18] A. Aghamohammadi et al., "A deep learning model for ergonomics risk assessment and sports and health monitoring in self-occluded images," Signal, Image and Video Processing, pp. 1-13, 2023.

[19] P. U. Patil, S. B. Lande, V. J. Nagalkar, S. B. Nikam, and G. Wakchaure, "Grading and sorting technique of dragon fruits using machine learning algorithms," Journal of Agriculture and Food Research, vol. 4, p. 100118, 2021.

[20] H. S. Gill and B. S. Khehra, "Fruit image classification using deep learning," 2022.

[21] M. Prem Kumar and A. Parkavi, "Quality grading of the fruits and vegetables using image processing techniques and machine learning: a review," Advances in Communication Systems and Networks: Select Proceedings of ComNet 2019, pp. 477-486, 2020.

[22] S. K. Chakraborty et al., "Development of an optimally designed real-time automatic citrus fruit grading–sorting machine leveraging computer vision-based adaptive deep learning model," Engineering Applications of Artificial Intelligence, vol. 120, p. 105826, 2023.

[23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770-778.

[24] N. M. Ibrahim, D. G. I. Gabr, A.-u. Rahman, S. Dash, and A. Nayyar, "A deep learning approach to intelligent fruit identification and family classification," Multimedia Tools and Applications, vol. 81, no. 19, pp. 27783-27798, 2022.

[25] S. Mukherjee. "The Annotated ResNet-50." https://towardsdatascience.com/the-annotated-resnet-50-a6c536034758 (accessed.

[26] M. Nikhitha, S. R. Sri, and B. U. Maheswari, "Fruit recognition and grade of disease detection using inception v3 model," in 2019 3rd International conference on electronics, communication and aerospace technology (ICECA), 2019: IEEE, pp. 1040-1043.

[27] L. Ali, F. Alnajjar, H. A. Jassmi, M. Gocho, W. Khan, and M. A. Serhani, "Performance evaluation of deep CNN-based crack detection and localization techniques for concrete structures," Sensors, vol. 21, no. 5, p. 1688, 2021.

[28] H. Altaheri, M. Alsulaiman, and G. Muhammad, "Date fruit classification for robotic harvesting in a natural environment using deep learning," IEEE Access, vol. 7, pp. 117115-117133, 2019.

[29] Y. Zheng, C. Yang, and A. Merkulov, "Breast cancer screening using convolutional neural network and follow-up digital mammography," in Computational Imaging III, 2018, vol. 10669: SPIE, p. 1066905.

[30] H. Herman, T. W. Cenggoro, A. Susanto, and B. Pardamean, "Deep Learning for Oil Palm Fruit Ripeness Classification with DenseNet," in 2021 International Conference on Information Management and Technology (ICIMTech), 2021, vol. 1: IEEE, pp. 116-119.

[31] N. Radwan, "Leveraging sparse and dense features for reliable state estimation in urban environments," University of Freiburg, Freiburg im Breisgau, Germany, 2019.

[32] S. Z. M. Zaki, M. A. Zulkifley, M. M. Stofa, N. A. M. Kamari, and N. A. Mohamed, "Classification of tomato leaf diseases using MobileNet v2," IAES International Journal of Artificial Intelligence, vol. 9, no. 2, p. 290, 2020.

[33] S. Phiphiphatphaisit and O. Surinta, "Food image classification with improved MobileNet architecture and data augmentation," in Proceedings of the 3rd International Conference on Information Science and Systems, 2020, pp. 51-56.

[34] N. Ismail and O. A. Malik, "Real-time visual inspection system for grading fruits using computer vision and deep learning techniques," Information Processing in Agriculture, vol. 9, no. 1, pp. 24-37, 2022.

[35] E. Cano, J. Mendoza-Avilés, M. Areiza, N. Guerra, J. L. Mendoza-Valdés, and C. A. Rovetto, "Multi skin lesions classification using fine-tuning and data-augmentation applying NASNet," PeerJ Computer Science, vol. 7, p. e371, 2021.

[36] T. Ahmed and N. H. N. Sabab, "Classification and understanding of cloud structures via satellite images with EfficientUNet," SN Computer Science, vol. 3, pp. 1-11, 2022.

[37] V. Bhole and A. Kumar, "Analysis of convolutional neural network using pre-trained squeezenet model for classification of thermal fruit images," in ICT for Competitive Strategies: CRC Press, 2020, pp. 759-768.

[38] E. Khan, M. Z. U. Rehman, F. Ahmed, and M. A. Khan, "Classification of diseases in citrus fruits using SqueezeNet," in 2021 International Conference on Applied and Engineering Mathematics (ICAEM), 2021: IEEE, pp. 67-72.

[39] Z. Guo, Q. Chen, G. Wu, Y. Xu, R. Shibasaki, and X. Shao, "Village building identification based on ensemble convolutional neural networks," Sensors, vol. 17, no. 11, p. 2487, 2017.

# Artificial Intelligence-driven Training and Improvement Methods for College Students' Line Dancing

Xiaohui WANG

Department of Arts and Sports, Henan Open University, Zhengzhou, Henan 450046, China

*Abstract*—With the advancement of computer technology, artificial intelligence technology has gradually become a research focus, and the thinking of relevant researchers has gradually transferred from the computer to the interaction between computers and humans. Artificial intelligence has begun to appear in various industries. With its rigorous computing logic and efficient computing speed, artificial intelligence technology begins to replace high-precision or highly repetitive work in work gradually. However, no specific data supports the specific work efficiency and output. In this context, this essay studies the methods of AI in the training and improvement of college students' line dancing levels. Virtual reality technology mainly undertakes functions such as virtual space modeling, sound positioning, sensory feedback, voice interaction, visual and spatial tracking, to ensure accurate positioning during choreography and motion capture. In this case, mechanical capture devices are used for motion capture in virtual reality space. This article uses intelligent capture technology based on virtual reality technology and artificial intelligence algorithms to capture and analyze the dance posture, generate analysis reports in a timely manner, and provide correction and suggestions for the dance posture. The final results show that AI can improve the training efficiency of line dancing of college students and can increase the innovation degree and method of dance posture by 7% to 13% compared with pure artificial. It shows that artificial intelligence technology plays a good role in college students' overall line dance training. At the same time, this paper also argues that artificial intelligence technology can effectively improve the overall productivity of traditional industries.

*Keywords*—*Motion capture; artificial intelligence technology; virtual reality; college students' line dance training; dance ascension*

## I. INTRODUCTION

Traditional row dance training usually adopts the method of teacher demonstration and student imitation, which may be difficult for some students with poor foundations to keep up and lack innovation and diversity. Traditional row dance training often focuses only on teaching dance steps, while neglecting training in music rhythm, body posture, dance expression, and other aspects, resulting in insufficient dance expression of students. In line dance training, teachers need to spend a lot of time demonstrating and correcting, and students also need to practice repeatedly to master, resulting in relatively low training efficiency. At present, the line dance training of college students in China mainly relies on traditional teaching methods, such as teaching in accordance

with textbooks or traditional chalk blackboards for pen presentations or wall charts, contrast training, and other methods. A common feature of these methods requires teachers or training leaders to demonstrate actions [1]. With the advancement of digital multimedia technology, the current teaching means gradually changing from traditional offline to online multimedia teaching. For example, common multimedia teaching means include video PPT display or other interactive ways of image and text. Multimedia teaching means can effectively save the energy of teachers or trainers. At the same time, the standard system of dance posture is established, which enriches teaching resources and plays a certain role in optimizing the whole teaching process [2]. But influenced by traditional teaching, many teachers do not like to use existing multimedia equipment in the course of line dance training. This is because multimedia equipment is a new teaching method compared with traditional teaching, and many teachers, especially older teachers, are unwilling to spend more energy or time preparing relevant learning materials. At the same time, there are some unfamiliar problems in the application and operation of the new equipment, which may lead to mistakes in the teaching process, thus reducing the teaching enthusiasm of teachers [3]. Schools or related personnel do not attach enough importance to multimedia teaching, coupled with the complexity of multimedia equipment, resulting in teachers' low enthusiasm for multimedia equipment teaching. On the other hand, when teachers use multimedia teaching, they replace the content of traditional blackboard writing with electronic information. At the same time, they still ignore students' enthusiasm in real classroom teaching or after-class training, resulting in students' passive acceptance of relevant knowledge. The relevant teaching technology based on artificial intelligence can effectively compensate for the above two points. First, the algorithm represented by artificial intelligence neural network can far surpass human eye recognition in the recognition and training standards of the accuracy of dance pose in the process of college students' line dancing and save a lot of manpower. At the same time, the technology represented by artificial intelligence virtual reality can transform two-dimensional classrooms into three-dimensional classrooms so that students can feel and learn from the scene and effectively improve their enthusiasm. Therefore, it is necessary to reform teaching methods from traditional dictation to two-dimensional blackboard writing and then to two-dimensional electronic information, and the transformation from two-dimensional classroom to

three-dimensional classroom can be called leapfrog progress, which is also the only way to reform teaching methods [4].

Artificial intelligence can provide personalized choreography training suggestions to students by analyzing their learning habits and levels, helping them better master dance techniques. Based on students' dance videos and action data, artificial intelligence can analyze their shortcomings and provide targeted training suggestions. Artificial intelligence can automatically analyze and evaluate students' dance videos, helping them quickly identify their shortcomings, thereby reducing repetitive practice time and improving training efficiency. Meanwhile, students can practice anytime and anywhere through smart devices, without being limited by time and location. Compared with China, foreign VR technology has been popularized and applied in all aspects of life, especially in the education industry. Virtual reality technology has been popularized in normal classroom education as a new multimedia technology in foreign classrooms. For example, Walmart in the U.S. used virtual reality technology to train 150,000 employees on basic professional operation requirements in 2017. They mainly use virtual reality technology to simulate a virtual trading scenario, then let employees simulate and practice it. India's Future Education 3.0 program already uses virtual reality to teach aviation courses. Based on the above information, the application of virtual reality is primarily seen in the education industry for simulating virtual scenes. With the simulation of a virtual scene, virtual reality technology can help the trainers train for the upcoming operation or scene earlier [5]. At the same time, virtual reality technology can also help us rehearse in advance the scenarios that cannot be simulated in practice, such as the simulation of aviation courses. By introducing AI technology into college dance training, we have the potential to overcome the limitations of traditional training methods and achieve a more personalized, efficient, and interesting training experience. However, we should also pay attention to and address the challenges and problems that come with it. Based on student feedback and performance, continuously optimize AI recommendation algorithms to improve the level of personalized guidance. Improve the real-time analysis and feedback capabilities of AI systems to ensure that students can promptly understand their own problems and make improvements. Combine the humanistic care of teachers, emotional communication, and objective data analysis of AI to form a more comprehensive and effective guidance method.

## II. RELATED WORK

At present, the research on the level of training and improvement methods of line dancing of college students based on AI technology is mainly based on two aspects: on the one hand, the new artificial intelligence technology represented by virtual reality, and on the other hand, the motion capture algorithm based on neural network. Virtual reality technology originated in the United States in the 1850s. It can be said that the Virtual reality technology in the United States represents the overall level of virtual reality technology in the world. Virtual reality technology was first applied in the military field. So far, VIRTUAL reality technology in the United States has been used in all aspects of life, but the United States pays particular attention to its application in education [6]. Zhang

and others based on virtual reality technology in the view of cognitive psychology to study the influence of the paper; he thinks that virtual reality technology teaching means through multiple senses for the instruction of learning knowledge, so enhances students' impressions when compared to traditional teaching methods, at the same time, cultivate the habit of students' autonomous learning and the search for knowledge and desire. Slatel et al. found through experimental studies that virtual reality technology can be closely combined with real sports to improve people's living standards while exercising [7]. Compared with the United States, the research on virtual reality technology in China was carried out later. Although the technology was introduced into China shortly after its birth, it did not get enough attention and attention at that time. However, in the early 1990s, the country and relevant researchers realized the importance of this technology. Since then, China's virtual reality technology has been developing in full swing [8]. Guo Xiaoming et al. first studied virtual reality glasses and ordinary glasses based on software Settings and analyzed the characteristics, functions, and structure of virtual reality glasses in detail, providing a theoretical basis for the subsequent generalization of virtual reality equipment. Song Da et al. introduced VR technology into education for the first time. He believed that virtual reality technology could serve as a new knowledge carrier, and students could reconstruct knowledge systems in an immersive state to cultivate students' ability for independent learning and exploration [9]. Overall, virtual reality technology research in China is still developing. On the one hand, virtual reality research requires expensive technical equipment, which limits the topic selection of relevant researchers regarding hardware. On the other hand, virtual reality is still a new technology in China, and most researchers are still skeptical about its role. Therefore, relevant researchers will deliberately avoid this topic in subjective topic selection. However, with the arrival of the information age, artificial intelligence represented by virtual reality has received more and more attention, and relevant research has gradually received strong support [10]. China's ninth Five-Year Plan and national Advanced Technology Development Strategy have repeatedly listed artificial intelligence technology represented by virtual reality as a key development project. As the frontier of scientific research, universities have gradually produced landmark achievements.

However, training line dancing level college students cannot play a practical role, only relying on virtual reality. The motion capture technology is also used to capture and analyze the relevant motion track to evaluate the actual motion and reverse output results, such as the accuracy of the motion, forming a training closed loop driven by artificial intelligence technology. Compared with virtual reality technology, motion capture technology develops later. It can be divided into motion capture methods, practical applications, and the analysis of motion data, posture, and other motion capture information [11]. Currently, motion capture technology combined with virtual reality technology has been gradually applied in intangible heritage protection, film and television production, game animation, medical rehabilitation engineering, and other aspects. For example, Qiu Wangbiao et al., based on virtual reality technology and motion capture algorithm, conducted data collection of different ethnic dances

and created relevant databases, which played a protective role in Chinese ethnic dance culture [12]. Kim et al. extracted individual walking characteristics based on a motion capture algorithm applied them to a humanoid biped robot, and analyzed their walking rules, step length, weight, and other data by recording them. Finally, the problem of walking friction individuals face in the rehabilitation process is effectively solved, and it can automatically recognize human walking posture and correct wrong posture in time [13]. Numerous studies have been conducted on motion capture technology in international sports education. For example, Wallance et al. carried out motion capture on the movement data of professional golfers in competitions and conducted professional analysis through the motion analysis system. Data captured in real-time through athletes' swing action, hitting posture, etc., provide theoretical support for subsequent athletes' further development and training direction [14]. Covaci et al. applied motion capture technology to volleyball training. Through data experimental analysis, they created a self-training shooting machine with virtual reality technology. It can not only analyze the posture of athletes, such as serving and receiving but also put forward optimization suggestions to athletes from angles and dynamics, which not only improves the training quality of volleyball class but also improves the training level of students [15]. Although artificial intelligence methods can provide personalized training suggestions and real-time feedback, they cannot replace the emotional communication and humanistic care of teachers. The words, deeds, encouragement, and care of teachers are of great significance for the growth and development of students. The update speed of artificial intelligence technology is very fast, and new algorithms and technologies are constantly emerging. In order to maintain the progressiveness and effectiveness of

technology, it is necessary to constantly upgrade and update the AI system, which requires a lot of human and material resources. Existing research mainly focuses on the feasibility of AI technology, while there is relatively little research on its practical effectiveness and long-term impact. Secondly, there is still insufficient research on how to combine the humanistic care of teachers with objective data analysis of AI to form more comprehensive and effective guidance methods. Finally, there is a lack of in-depth exploration on the acceptance and response of different students to AI assisted training.

## III. METHOD

To sum up, there are few application cases of virtual reality technology for classroom teaching with high accuracy. For instance, virtual reality technology is limited to VR instructional videos. It is speculated that this is because, for the training of related sports movements, the virtual reality technology needs to capture and analyze the relevant movements in the virtual scene to improve the real training level. However, the movement of the human body is a complex system involving the precise movement of multiple nodes or muscles. The current development of artificial intelligence algorithms is not enough to accurately and reasonably analyze human movement. However, suppose artificial intelligence technology based on virtual reality can be applied in sports training. In that case, it can not only improve the quality of physical education but also give students an immersive experience to ensure students' enthusiasm and training results in the process of training. Therefore, this paper takes the training and improvement methods of line dancing level of college students as an illustration to carry out pertinent artificial intelligence research represented by virtual reality. The research ideas of this paper are shown in Fig. 1.



Fig. 1. Research idea based on artificial intelligence technology driving line dance training of college students.

From Fig. 1, it can be seen the key to improving the level of line dance training of college students, which is also the data basis of this study, is to capture information on human posture during line dance training of college students through motion capture technology. Because multiple dancers are involved in the training process of line dancing, each dancer's movements may be different. Still, each slight movement may affect the overall performance of line dancing, so the requirements for movement capture are more accurate. Human body posture is mainly completed through joint muscles and other common performances. Each part has a general movement posture, but there are slight differences because dance movements for the whole body coordination requirements are relatively high. In the past, there was mainly attitude analysis based on simulated annealing particle swarm optimization, model-based motion attitude analysis, and feature-based motion attitude analysis for the human body. The pose analysis based on the simulated annealing particle swarm optimization algorithm is mainly for motion capture of human arm movement; a particle-wave algorithm is used to measure the conditional density of sampling points, and the random weight parameters are updated using the Monte Carlo method. Finally, the arm movement state of the human body is judged by the random particle movement. The specific algorithm is shown in Formula 1.

$$E(y_t, x_t) = S \times \left( (1 - \beta) \frac{B_t}{B_t + Y_t} + \beta \frac{R_t}{R_t + Y_t} \right) \qquad (1)$$

The capture of human body pose by model-based motion attitude analysis is mainly by establishing a model in advance, mapping the standard parameter model to the actual captured trajectory of human body pose, and calculating the difference between them by function to determine the deformation degree of the two to judge the human body pose. Therefore, the accuracy of model-based motion attitude analysis results depends more on the model's standard degree and the function's fitting effect. The fitting method is usually the least square method, and the polynomial coefficients of the curve fitting are used to approximately represent the rules of human movement. The specific functions are shown in Formula 2.

$$s(t) = a_0 \varphi_0 t + a_1 \varphi_1 t + \cdots + a_n \varphi_n t \qquad (2)$$

The feature-based motion pose analysis algorithm does not need to build a model as a whole but can directly determine the position of the human body's changing pose through the comparison of human pose features. For example, points, lines, planes, nodes, and other more complex features are used for comparison and reference.

The line dance training of college students to be explored in this paper involves many individuals, and the feature reference is extremely complex. If the feature motion pose analysis algorithm is used alone, the accuracy of pose change cannot be guaranteed. However, if the annealing particle swarm optimization algorithm is used to predict each individual's changing characteristics in the line dancing process. In that case, many operations will be produced, and the efficiency of computer operations will be seriously affected. Therefore, after comprehensively considering the balance between accuracy and efficiency, the model-based motion pose analysis

algorithm can be chosen as this paper's main motion capture algorithm. Nevertheless, the function fitting algorithm utilizes the least square method, requiring us to identify the optimal function parameters within the system space expectation to guarantee the precision of the fitting. The expression of spatial expectation is shown in Formula 3.

$$\phi = span\{\varphi_0 t, \varphi_1 t, \dots, \varphi_n t\} \qquad (3)$$

However, in the actual least square method fitting process, It was discovered that the sum of squares generated by the fitting increased as the number of iterations increased, indicating that the accuracy of the fitting became lower and lower as the number of iterations increased, as shown in Fig. 2. It is speculated that this is because the complexity of dance movements leads to more and more errors in the traditional least square fitting. Therefore, the traditional least square fitting method cannot be applied to the line dance training of college students studied in this paper. More precise fitting methods are required.

Therefore, an artificial intelligence-based neural network algorithm is introduced here. Currently, commonly used neural networks are mainly separated into two groups -- convolutional neural networks (CNN) and cyclic neural networks. The basic algorithm structure of the two algorithms is the same. The algorithm structure is divided into an input layer, an output layer, and a hidden layer. The input and output layers are the algorithm layers of initial data input and final result output at the beginning of the operation. However, the two have different logic in the hidden layer. CNN does not participate in the cycle during the training of the hidden layer, so the concept of timing is not considered, and the algorithm does not have the logical ability to relate to the context. However, the cyclic neural network will carry out cyclic training simultaneously in the hidden layer algorithm training so that the cyclic neural network will consider the time sequence in the training process. Normally, the convolutional neural network is used for image or image recognition because the input value is directly related to the output value in the recognition algorithm, and the influence of the output result of the previous neuron is not considered when the output neuron is generated in the next result. However, the recurrent neural network is usually used to generate natural language because it considers the temporal influence, that is, the influence of the output result of the last neuron. Therefore, the convolutional neural network is mainly introduced in this paper. In the traditional least square fitting process, the error value increases as the number of iterations increases, which is called a gradient explosion problem in the realm of AI and corresponds to a gradient attenuation problem. Both refer to the problem that the original algorithm is invalid and the fitting is out of order because the operation exceeds the threshold value. The artificial intelligence concept of a control gate is introduced to solve this problem. Specifically, it refers to adding two propagation control domain restrictions in artificial intelligence's hidden and input layers, called the update and reset gates. The specific expressions are shown in Formula 4 and 5.

$$r_t = \sigma(W_r * [h_{t-1}, x_t]) \qquad (4)$$

$$z_t = \sigma(W_z * [h_{t-1}, x_t]) \qquad (5)$$

Fig. 2. The traditional least square method is used to fit the curve.

The function of the update gate is to control the proportion of information output from the previous neuron to the current neuron. When the value of the update gate is set larger, it means that the output information of the previous neuron is input to the current neuron at a larger proportion. The function of the reset gate is to ignore the information input to the current neuron again. The smaller the reset gate value is, the more the current input content is ignored. The overall control gate follows the law of forward propagation. At the same time, the training network's loss function must be computed after passing the control gate each time, which also corresponds to the many-to-many or one-to-one model of a convolutional neural network. See Formulas 6 to 8 for specific propagation modes. The fitting results before and after the control was added are shown in Fig. 3.

$$h'_t = tanh(W_{h'} * [r_1 * h_{t-1}, x_t]) \qquad (6)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * h' \qquad (7)$$

$$y_t = \sigma(W_o * h_t) \qquad (8)$$



Fig. 3. The results of the least square method fitting are introduced.

It can be seen from Fig. 3 that before adjustment, the sum of squares in different series of fitting processes increased with the increase of fitting series. In contrast, after control gate adjustment, the fitted sum of squares showed an increasing and decreasing trend. This is considered because the content of acceptable information is large at the beginning due to the computer's large amount of free memory, so the control gate does not play a role for the time being. The square deviation of

the fitting keeps increasing before filtering out some invalid information. However, as the number of fitting iterations increases, the computing load of the computer becomes larger and larger. At this time, the control gate starts to operate and continuously reduces the input information of the previous neuron by updating the gate. Meanwhile, the invalid input information is ignored by resetting the gate to control the volume and effectiveness of the fitting information. Finally, the square deviation of the least square method is reduced, and the gradient explosion problem is avoided effectively.

After determining the fitting method, the algorithm of human pose recognition, the core of motion capture began to improve. The traditional pose recognition algorithm mainly refers to the 3d model similarity matching algorithm to measure the pose difference or similarity between different human bodies. It first carries on the special point identification to the human posture to be recognized. It calculates the change difference before and after the identification point through the Euclidean distance when it changes during the human body's line dance training to preliminarily determine the motion trajectory. The calculation method is shown in Formulas 9 to 10.

$$D = sqrt((x_1 - x_2)^2 + (y_1 - y_2)^2) \qquad (9)$$

$$D = sqrt((x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2) \quad (10)$$

The principle of similarity matching algorithm for different 3D models based on Euclidean distance is to compare two object representation points on the motion trajectory, get the corresponding distance difference through comparison, and get the corresponding human body posture by matching with the pre-set track-attitude comparison table. However, this algorithm has some shortcomings. First, identifying the point comparing method due to the human body posture is usually completed jointly by multiple joints and muscles, so to ensure the accuracy of the attitude to capture, it needs to set up multiple identification points to action and compare multiple points to identify at the same time, the amount of calculation is too big. Secondly, there is a certain stubbornness in determining attitude according to the scale method. In the line dance training of college students, dancers of different heights and weights have different changes in the identification points on their bodies when performing the same movement demonstration. However, there is a certain error in attitude comparison and confirmation according to the same scale at this time. The comparison of the movement track of the same identification point in a unilateral upward direction is shown in Fig. 4.



Fig. 4. A single identity between unilateral upward motion path.

It can be seen from Fig. 4 that the single-direction motion trajectory between single marker points is normally distributed over time. It indicates that with the increase of time, the Euclidean distance begins to lose regularity in the identification of marker locus, and only the irregular discrete marker locus changes with time will present normal distribution. In line dancing of college students, more individual pose recognition will be involved. Compared with the analysis of individual pose recognition, the computation amount of this method increases exponentially. Therefore, finally, abandoning the traditional algorithm and adopting the similarity calculation algorithm of feature plane matching has been decided.

The feature plane matching algorithm first extracts the skeleton model of the human body by motion capture technology, takes the extracted skeleton model as the basic computing plane, and marks all important parts of the human body, such as joints. The feature plane matching algorithm

differs from the traditional Euclidean distance algorithm in that the individual feature plane normal vector judges the attitude change. At the same time, the Angle between Eigen plane edge vectors is introduced to further judge the local normalization of motion attitude. The specific calculation method is from Formula 11 to Formula 13. According to this algorithm, the standard included Angle parameters of each comparison point of human joints can be obtained, as shown in Fig. 5.

$$similarity(V_i, V_j) = \frac{v_i \times v_j}{\sqrt{v_i^2} \times \sqrt{(v_j)^2}} \qquad (11)$$

$$\theta_{\langle i,j \rangle} = arccos\left(similarity(V_i, V_j)\right) \qquad (12)$$

$$Corr(\theta_{\langle i,j \rangle}) = 1 - \left(\frac{arccos(similarity(V_i,V_j))}{\pi}\right) \qquad (13)$$

Fig. 5. Each joint of the human body corresponds to the bending angle.

The matching results of each joint's bending Angle and plane feature show that there are still some errors between the predicted difference and the actual difference of some joints. These errors are the fundamental cause of inaccurate attitude analysis. After specific analysis, it was found that in the process of plane Angle comparison, the included Angle relation between limbs and trunk in the vertical direction was not taken into account, which led to the error in the mapping of the local relation between limbs and trunk, which led to the phenomenon of error increase. Therefore, a binary group is introduced here as a calculation model for locally fitting the similarity between limbs and trunk. Through the determination and output of relevant parameters, errors caused by the inherent characteristics of the human body are effectively solved. The specific expressions are shown in Formula 14 to Formula 15. Fig. 6 displays the outcome. The algorithm suggested in this paper reduces computational complexity and improves the efficiency and improvement method of AI-based university students' line dancing level training.

$$R^{(uv)} = \begin{bmatrix} R_{b_w, b_v} & \cdots & R_{b_w, b_{v+1}} \\ \vdots & \ddots & \vdots \\ R_{bu+1, b_v} & \cdots & R_{bu, b_{v+1}} \end{bmatrix} \tag{14}$$

$$S_{u,v} = \frac{Q_{max}^{(uv)}}{min(l_u, l_v)} \tag{15}$$



Fig. 6. The actual and predicted values of different mapping angles are displayed.

## IV. RESULT ANALYSIS

In our research, we found that significantly improved row dance performances are mainly manifested in the following aspects:

Accuracy and fluency of actions: Students who have undergone AI assisted training have significantly improved their accuracy and fluency of actions. The AI system can accurately identify subtle deviations in movements and provide immediate feedback and correction suggestions by analyzing

students' dance videos in real-time. This enables students to correct errors in a timely manner during the training process, gradually improving the accuracy and fluency of their movements.

Rhythm and music coordination: AI assisted training not only focuses on the accuracy and skills of movements, but also focuses on cultivating students' sense of rhythm and music coordination ability. Through precise analysis of music rhythm through AI systems, students can better understand the coordination between dance and music, and improve their performance in rhythm changes and complex rhythm types.

Creativity and expressiveness: AI technology not only provides traditional training methods and guidance, but also stimulates students' creativity and expressiveness through data analysis. For example, AI systems can analyze students' dance movements, generate personalized dance choreography suggestions, and help students explore different dance styles and forms of expression.

The fundamental procedure and algorithm of artificial intelligence technology powering the level training and enhancement approach of line dancing for college students have been validated and tested above. Next, the algorithm will be verified through the practical operation and the output of the

final results. The need to introduce virtual reality devices; virtual reality technology here mainly undertakes virtual space modeling, sound localization, sensory feedback, voice interaction, visual and space tracking, and other functions to ensure accurate positioning in the process of line dance and motion capture, in this instance, mechanical capture devices are used for motion capture in virtual reality space. The captured data is analyzed for action, as shown in Fig. 7.

To ensure the objectivity and effectiveness of the experiment, training personnel were randomly divided into two parts before the beginning of the line dance training. One part still adopts the traditional line dance training method; the teacher teaches offline, and the students do not use any artificial intelligence equipment or methods for training. Others use artificial intelligence to train for line dancing. On the one hand, the virtual reality equipment is used for unlimited time and place training in the virtual scene; on the other hand, the intelligent capture technology based on an AI algorithm is used to capture and analyze the line dance posture and timely produce analysis reports for correction and suggestions on dance posture. After one month of such training, each student's dance posture is evaluated again - and the evaluation results are shown in Fig. 8.



Fig. 7. An example of motion capture in a virtual reality scene based on artificial intelligence technology.



Fig. 8. The score of the same trainer before and after the training using artificial intelligence.

First, the level of the students who used artificial intelligence technology for line dancing training before and after the training has been compared. Six teachers have been invited as judges to conduct a manual evaluation and score the dancing posture of the students who used artificial intelligence technology for line dancing training before and after the training. It can be seen from Fig. 8 that the scores of students before and after training are generally higher than those before training. Through calculation, it can roughly be concluded that the line dancing level of the students who use artificial intelligence technology for line dancing training has improved by about 20% on average. However, this result does not mean that adopting AI technology causes all students to improve their training. Even without external forces' help, there is only an improvement effect through their practice. Therefore, the level test of students before and after training was conducted to eliminate the impact of individual objective elements without

the aid of AI equipment, and the results showed that the level of these students only improved by about 10%, which, in our opinion, was the increase of individual natural level without the help of equipment. Therefore, artificial intelligence technology has an effect of 7%~13% on improving the level of line dancing training of college students.

Finally, to get everyone to the proposed artificial intelligence technology drive line dance training level college students to improve the algorithm's true feelings, whether to accept the technology popularization and application in the classroom and the benefits of the technology are two basic problems for everyone, the results showed the 51 people welcome the introduction of the technology in education, It is considered that this technology plays a positive role in promoting students' innovation and autonomy, as shown in Fig. 9.



Fig. 9. Artificial intelligence technology drives acoustic feedback of line dance training for college students.

## V. DISCUSSION

The accuracy of body posture and dance movements is a key factor in the quality of performance. Through artificial intelligence methods, we can accurately analyze student movement data and identify subtle but key improvements such as posture stability, movement fluency, and coordination. This specific feedback information can help students more accurately understand their shortcomings and make targeted improvements.

Secondly, the sense of musical rhythm is an important component of line dance performance. Artificial intelligence can provide feedback on rhythm control and coordination between dance and music by analyzing music rhythms and student dance movements. This helps students improve their understanding and expressive power of music, further enhancing the overall dance effect.

As for the analysis of variability, preliminary studies have shown that there is a certain degree of variability in the

reactions of different students to artificial intelligence assisted training. Some students have shown a high degree of acceptance and active cooperation towards new technologies, believing that artificial intelligence provides personalized and efficient training methods; some students, on the other hand, rely more on traditional training methods and hold a reserved attitude towards artificial intelligence. This variability may be related to students' technological acceptance, learning style, and habits.

## VI. CONCLUSION

The research mentioned above reveals that artificial intelligence technology is being used in education but is not popular. Still, in any sports teaching or training research, it is inseparable from the motion capture system to record and examine the posture of the human body if the scene is needed simultaneously with the help of virtual reality technology to build the virtual scene. So, this article for the college student's level of line dance training and the promotion of the research is based on the idea of method, first through the motion capture

system in the process of line dance teaching and training for college students in data capture and analysis, according to the certain algorithm to calculate its standard and accuracy, at the same time on the place of the action is not standard tag and automatically in the course of the next training initiative to remind. In this way, the level of line dancing training can be improved. In the virtual scene constructed by VR technology, the innovation of a new dance pose can be simulated at will, and the effect can be watched synchronously. Therefore, it is theoretically expected that the improvement of training levels and innovation of training methods and systems regarding the impact of AI technology on the line dance training of college students. The final results show that AI can improve the training efficiency of line dancing of college students and can increase the innovation degree and method of dance posture by 7% to 13% compared with pure artificial.

## COMPETING OF INTERESTS

The authors declare no competing of interests.

## AUTHORSHIP CONTRIBUTION STATEMENT

Xiaohui Wang: Writing-Original draft preparation, Conceptualization, Supervision, Project administration.

## REFERENCES

[1] P. Aimsamrarn, T. Janyachareon, K. Rattanathanthong, A. Emasithi, and W. Siritaratiwat, "Cultural translation and adaptation of the Alberta Infant Motor Scale Thai version," Early Hum Dev, vol. 130, pp. 65–70, 2019.

[2] M. Yang, S. Liu, K. Chen, H. Zhang, E. Zhao, and T. Zhao, "A hierarchical clustering approach to fuzzy semantic representation of rare words in neural machine translation," IEEE Transactions on Fuzzy Systems, vol. 28, no. 5, pp. 992–1002, 2020.

[3] B. Yang, D. F. Wong, L. S. Chao, and M. Zhang, "Improving tree-based neural machine translation with dynamic lexicalized dependency encoding," Knowl Based Syst, vol. 188, p. 105042, 2020.

[4] Q. Gao and J. Zhou, Human Aspects of IT for the Aged Population. Technology in Everyday Living: 8th International Conference, ITAP 2022, Held as Part of the 24th HCI International Conference, HCII 2022, Virtual Event, June 26–July 1, 2022, Proceedings, Part II, vol. 13331. Springer Nature, 2022.

[5] W. Ho, Culture, Creativity, and Music Education in China: Developments and Challenges. Taylor & Francis, 2023.

[6] H. Zhang, "Research on dancer tracking technology based on contour model and AdaBoost algorithm," in International Conference on Mechanisms and Robotics (ICMAR 2022), SPIE, 2022, pp. 1376–1382.

[7] C. Jie and T. Feng, "Design and implementation of virtual dance training system based on Unity3D [J]," Industrial Control Computer, vol. 32, no. 012, pp. 49–51, 2019.

[8] Y. Mirsky and W. Lee, "The creation and detection of deepfakes: A survey," ACM Computing Surveys (CSUR), vol. 54, no. 1, pp. 1–41, 2021.

[9] R. D. S. Yates and D. Cai, "Bibliography of Women and Gender in China (2018-2022)," NAN NÜ, vol. 25, no. 2, pp. 213–343, 2023.

[10] L. Huang, "Application status and design requirements of virtual reality technology in dance teaching," in 2019 International Conference on Advanced Education, Service and Management, The Academy of Engineering and Education, 2019, pp. 730–733.

[11] W. Fu, S. Liu, and J. Dai, "E-learning, e-education, and online training," 2018.

[12] X. Lu, "A Feasibility Study of VR Technical Practice in Dance Teaching [J]," Art Education, vol. 14, pp. 99–100, 2018.

[13] H. Zhang, "Research on dancer tracking technology based on contour model and AdaBoost algorithm," in International Conference on Mechanisms and Robotics (ICMAR 2022), SPIE, 2022, pp. 1376–1382.

[14] M. Skublewska-Paszkowska, M. Milosz, P. Powroznik, and E. Lukasik, "3D technologies for intangible cultural heritage preservation—literature review for selected databases," Herit Sci, vol. 10, no. 1, pp. 1–24, 2022.

[15] Y. Mirsky and W. Lee, "The creation and detection of deepfakes: A survey," ACM Computing Surveys (CSUR), vol. 54, no. 1, pp. 1–41, 202.

# State-of-the-Art Review of Deep Learning Methods in Fake Banknote Recognition Problem

Ualikhan Sadyk, Rashid Baimukashev, Cemil Turan

Suleyman Demirel University, Kaskelen, Kazakhstan

*Abstract*—In the burgeoning epoch of digital finance, the exigency for fortified monetary transactions is paramount, underscoring the need for advanced counterfeit deterrence methodologies. The research paper provides an exhaustive analysis, delving into the profundities of employing sophisticated deep learning (DL) paradigms in the battle against fiscal fraudulence through fake banknote detection. This comprehensive review juxtaposes the traditional machine learning approaches with the avant-garde DL techniques, accentuating the conspicuous superiority of the latter in terms of accuracy, efficiency, and the diminution of human oversight. Spanning multiple continents and currencies, the discourse highlights the universal applicability and potency of DL, incorporating convolutional neural networks (CNNs), recurrent neural networks (RNNs), and generative adversarial networks (GANs) in discerning the most cryptic of counterfeits, a feat unachievable by obsolete technologies. The paper meticulously dissects the architectures, learning processes, and operational facets of these systems, offering insights into their convolutional strata, pooling heuristics, backpropagation, and loss minimization algorithms, alluding to their consequential roles in feature extraction and intricate pattern recognition - the quintessentials of authenticating banknotes. Furthermore, the exploration broaches the ethical and privacy concerns stemming from DL, including data bias and over-reliance on technology, suggesting the harmonization of algorithmic advancements with robust legislative frameworks. Conclusively, this seminal review posits that while DL techniques herald a revolutionary competence in fake banknote recognition, continuous research, and multi-faceted strategies are imperative in adapting to the ever-evolving chicanery of counterfeit malefactors.

*Keywords—Fake banknote; detection; classification; recognition; review*

## I. INTRODUCTION

Counterfeiting remains one of the most insidious challenges facing monetary institutions worldwide, with its implications stretching beyond mere economic effects to encompass significant social and security dimensions. The global prevalence of counterfeit currency has witnessed an alarming increase, with the Financial Action Task Force (FATF) and the International Monetary Fund (IMF) highlighting the substantial threats posed by this illicit activity to the integrity of financial markets and, by extension, national security [1]. The sophistication of modern counterfeiting techniques, enabled by technological advancements, necessitates an equally advanced approach to counterfeit currency detection and prevention.

Traditional methods of counterfeit detection have revolved around manual and mechanical authentication techniques, ranging from the scrutiny of security features visible to the naked eye to the use of rudimentary electronic validators. These methods, although somewhat effective in the past, are increasingly falling short in the face of advanced counterfeiting. Studies indicate that conventional methodologies demonstrate limited success, especially with the advent of high-definition color printing and the replication of primary security features, often failing to catch more sophisticated counterfeit notes and leading to a significant volume of false negatives [2].

Moreover, the human factor in traditional methods often results in inconsistencies; studies have revealed that repetitive tasks combined with high-pressure environments significantly increase human error, leading to lapses in detection [3]. Similarly, mechanical validators are constrained by their programming based on specific features of banknotes. They do not adapt to new security enhancements without reprogramming or replacement, making them both economically and operationally inefficient in the long run [4].

In contrast, the emergence of deep learning techniques has heralded a transformative approach to counterfeit detection. Deep learning, a subset of machine learning, is characterized by algorithms that mimic the neural circuitry of the human brain to progressively improve performance on tasks [5]. Within the sphere of counterfeit detection, deep learning models, particularly Convolutional Neural Networks (CNNs), have demonstrated the capability to identify subtle inconsistencies and deviations on banknotes, which would typically go unnoticed by human inspectors or conventional machinery [6].

One of the most significant advantages of integrating deep learning into counterfeit detection is its ability to learn and adapt continually. These systems are designed to evolve with every data point, enhancing their accuracy over time and allowing them to keep pace with emerging counterfeiting technologies without the need for manual intervention or reprogramming [7]. Additionally, they reduce the cognitive load and error rate associated with human inspection, thereby streamlining the verification process [8].

However, the application of deep learning is not without challenges. The efficacy of these systems is heavily reliant on the availability and quality of training data, necessitating extensive datasets of both counterfeit and genuine banknotes for initial setup and ongoing learning [9]. Despite these requirements, the potential of deep learning in revolutionizing

banknote authentication practices is gaining recognition, with several central banking institutions and financial bodies investing in this technology [10].

This review paper aims to provide a comprehensive overview of the application of deep learning techniques in the detection of counterfeit banknotes. It seeks to explore the evolution from traditional methods to advanced technological means, emphasizing the increasing inadequacy of the former and the promising capabilities of the latter. The review will delve into various deep learning models, examining their operational mechanisms, advantages, and potential limitations in the context of counterfeit detection [11].

Furthermore, this paper will analyze real-world applications and case studies where deep learning techniques have been successfully implemented. It will highlight the practical considerations and logistical implications of integrating these systems into existing financial security frameworks [12]. In doing so, it will also touch upon the challenges, particularly those related to ethics and data security, that come with the adoption of advanced AI technologies in sensitive sectors. Fig. 1 demonstrates a sample of fake banknote detection system [13].

By drawing upon a wide range of sources, including scholarly articles, industry reports, and white papers [14-18], this review intends to offer a multi-dimensional perspective on the subject. It is directed towards academics, professionals, and decision-makers in the fields of finance, security, and artificial intelligence, providing them with a consolidated resource that not only underscores the urgency of adopting more sophisticated counterfeit detection methods but also guides future research and policy-making in this critical domain.



Fig. 1.   Fake banknote detection system.

## II.   TRADITIONAL METHODS FOR COUNTERFEIT DETECTION

The historical landscape of combating monetary forgery has primarily relied on several traditional methods, each with distinct mechanisms designed to discern the authenticity of banknotes. These conventional strategies, while having served financial institutions for decades, exhibit certain limitations, especially in the face of technologically advanced counterfeiting tactics [19].

One of the most longstanding techniques is the use of watermark technology, where an image or pattern is embedded into the physical structure of the paper itself. This method, requiring the transmittance of light through the note for verification, has been a hallmark of banknote security. However, with advancements in digital imaging and printing technology, counterfeiters have been able to simulate watermarks to a convincing degree, diminishing the effectiveness of this once-reliable method [20]. Fig. 2 demonstrates flowchart of an image processing for counterfeit detection system [13].



Fig. 2.   Sample flowchart of a counterfeit detection system.

Security threads integrated into the substrate of banknotes comprise another traditional safeguard against counterfeiting. These metallic or plastic threads, often partially embedded and partially exposed, are designed to be distinctive and challenging to replicate. Despite this, modern counterfeiting operations, utilizing advanced materials and manufacturing techniques, have successfully imitated such features, leading to the circulation of fake notes undetected by standard thread verification processes [21].

Ultraviolet (UV) features, visible only under UV light, and micro-printing, where minute text or images are printed on the banknote, have also been employed historically. While these features are less accessible for replication by amateur

counterfeiters, organized and technologically equipped counterfeit operations have managed to bypass these security measures. The mass production of counterfeit notes with passable UV features and micro-printing has exposed the vulnerability of these methods [22].

Additionally, the feel of the paper, raised printing, and other tactile elements have long been the first line of defense, as cash handlers traditionally rely on touch to detect counterfeit notes instinctively. The reliance on sensory perception, albeit practical and cost-effective, is highly subjective and prone to human error. The introduction of high-grade counterfeit notes, mimicking the tactile features of genuine banknotes, complicates the reliability of this sensory approach [23].

The use of magnetic ink and the magnetic properties of certain printing elements present on genuine banknotes has been a cornerstone of automated banknote validation within vending machines and note counters. Counterfeiters have, however, found ways around this through the application of magnetic ink in appropriate areas, confusing the sensors and limiting the success of magnetic detection [24].

Moreover, traditional methods face a common limitation: the need for human intervention, whether in the direct handling and inspection of notes or in the maintenance and updating of machinery used for detection. This human dependency increases the likelihood of inconsistency and error, thereby reducing the overall efficacy of counterfeit detection measures [25].

The advancements in counterfeiting technology, alongside the limitations of traditional detection methods, highlight an arms race between counterfeiters and authorities. As counterfeiters adopt more sophisticated technology, they exploit the weaknesses inherent in traditional methods, necessitating a move towards more advanced, technology-driven detection systems [26].

In light of these insights, financial institutions and governing bodies have been impelled to explore and adopt more technologically advanced methods, particularly in the realm of artificial intelligence and machine learning. The transition is driven by the need to enhance accuracy, speed, and adaptability in the detection processes—attributes that are increasingly pertinent in the context of modern, sophisticated counterfeiting techniques [27].

Conclusively, while traditional methods have played a significant role in counterfeit detection, their relevance is waning in the current technological climate. The limitations and challenges they present underscore the necessity for innovation and advancement in this field, pointing towards deep learning and other AI methodologies as the next logical step in counterfeit detection [28-32]. This transition is not just a matter of enhancing efficiency, but an imperative adaptation for maintaining the integrity of global currency systems in the contemporary age.

III. EMERGENCE OF DEEP LEARNING IN COUNTERFEIT DETECTION

The relentless evolution of counterfeiting practices has necessitated a paradigm shift in detection methodologies, steering the discourse towards more resilient, adaptable, and sophisticated solutions. At the forefront of this evolution is deep learning, a revolutionary approach that has transcended the theoretical boundaries of computer science to establish itself as an instrumental asset in practical counterfeit detection.

A. *Definition and General Concept of Deep Learning*

Deep learning, a subset of machine learning in artificial intelligence (AI), orchestrates learning from data that is unstructured or unlabeled at colossal scales. It employs algorithms operating in layered structures known as neural networks, which are designed to imitate the human brain's decision-making process [33]. Each layer of a neural network filters inputs from expansive datasets, making independent decisions on the data and passing it to the next layer. Through this hierarchical processing, deep learning models can make sense of large-scale data with complex patterns, a feat unattainable by traditional machine learning models. Fig. 3 demonstrates a sample of counterfeit detection system process using deep learning technologies [34].

Unlike standard machine learning models that plateau in performance as more data is supplied, deep learning models continue to improve. This characteristic is crucial in scenarios where data is abundant, and subtle nuances in data are vital for making accurate predictions or classifications, such as distinguishing genuine banknotes from counterfeits [35].



Fig. 3. Sample flowchart of a counterfeit detection system using deep learning [34].

## B. Historical Context in the Field of Artificial Intelligence (AI)

The conception of deep learning dates back to the 1940s, with the advent of the "perceptron" — the simplest form of a neural network, capable of learning and making decisions on its own [36]. However, it was not until the 1980s that interest in neural networks resurged, attributed to the backpropagation algorithm, which allowed networks to adjust hidden layers of neurons in an efficient manner [37].

Despite these advancements, early neural networks were rudimentary, with their learning capabilities limited by the computational power and data availability of the time. The dawn of the 21st century, marked by a digital explosion and unprecedented advancements in computational power, set the stage for today's deep learning landscape. This era witnessed the convergence of a massive influx of data (big data) and significantly enhanced computing capacities, including the use of Graphics Processing Units (GPUs) to fast-track deep learning computations [38].

## C. Emergence of Deep Learning in Counterfeit Detection.

Healthcare. In healthcare, deep learning has been a catalyst for innovation, particularly in medical imaging. Deep learning models, through pattern recognition, have significantly improved the diagnosis, prognosis, and treatment planning of diseases, matching, and occasionally surpassing expert-level accuracy [39]. For instance, convolutional neural networks (CNNs) have demonstrated remarkable precision in detecting skin cancer, diabetic eye diseases, and other pathologies from medical images, underscoring the potential of deep learning in enhancing medical diagnostics [40].

Autonomous Vehicles. The autonomous vehicle industry has leveraged deep learning to improve navigation and safety. By processing vast datasets from various sensors and cameras, deep learning systems can make split-second decisions on the road, recognizing objects, predicting pedestrian movements, and identifying potential hazards. This continuous learning process is pivotal for the development of safe, reliable autonomous vehicles [41].

Finance. The finance sector, characterized by its dynamic and complex nature, has employed deep learning for various applications including algorithmic trading, risk management, and customer service. Neural networks process market indicators efficiently, providing insights for investment and trading decisions [42]. Furthermore, AI-driven chatbots, powered by deep learning, handle customer inquiries, process transactions, and detect fraudulent activities, offering enhanced efficiency and security [43].

Cybersecurity. Deep learning's application in cybersecurity has transformed threat detection by analyzing network traffic, identifying unusual patterns, and mitigating threats in real-time. Traditional cybersecurity measures struggle to keep pace with the sophistication of modern cyber-attacks, but deep learning models thrive on this complexity, continually adapting and learning from new data [44].

Retail. The retail sector harnesses deep learning for personalized shopping experiences, inventory management, and logistics. AI models analyze customer data, predicting shopping trends, and behavior to recommend products uniquely suited to individual preferences, significantly driving sales and customer satisfaction [45].

Manufacturing. In manufacturing, deep learning facilitates predictive maintenance, quality control, and supply chain optimization. By predicting machine failures before they occur, companies can plan maintenance without disrupting production, a testament to deep learning's preventative potential [46].

These diverse applications underscore deep learning's adaptability and its transformative impact across industries. Its ability to decipher complex patterns from vast datasets, predict outcomes, and automate decision-making processes is universally beneficial. As counterfeit detection techniques integrate deep learning, they leverage these strengths, offering improved accuracy, adaptability, and reliability in distinguishing genuine banknotes from sophisticated forgeries [47-58]. The versatility of deep learning, evidenced by its broad utilization, not only enhances the capabilities within each respective field but also contributes profoundly to the advancement of interdisciplinary technological innovations.

## IV. DEEP LEARNING TECHNIQUES FOR FAKE BANKNOTE RECOGNITION

The burgeoning field of deep learning has ushered in innovative techniques that significantly enhance the accuracy and efficiency of counterfeit banknote recognition. These methodologies, grounded in different aspects of artificial intelligence, have been pivotal in revolutionizing the approach towards ensuring the authenticity of currency.

### A. Convolutional Neural Networks (CNNs)

Structure and Functionality. Convolutional Neural Networks (CNNs) are a class of deep neural networks that have become the gold standard for image recognition tasks, owing to their architecture optimized for processing grid-like data, including pixels in images [59]. CNNs consist of multiple layers, notably the convolutional layers, which use filters to create feature maps that retain spatial relationships across the input, capturing the dependencies among pixels in close proximity. These layers are complemented by pooling layers, reducing computational load, and controlling overfitting by progressively downsizing the spatial dimensions of the input representation [60].

Suitability for Image Recognition. CNNs stand out in image classification and object detection due to their ability to automate feature extraction from raw data, a process that traditional algorithms could not perform without extensive manual feature engineering [61]. When applied to banknote verification, CNNs can analyze intricate details in banknotes, discerning genuine features from counterfeit attempts by learning discriminative features, which are often overlooked by the human eye and traditional computational methods [62].

### B. Recurrent Neural Networks (RNNs)

Operational Mechanism. Recurrent Neural Networks (RNNs) are another subset of neural networks where connections between neurons form a directed graph along a

sequence, allowing it to exhibit temporal dynamic behavior. Unlike traditional neural networks, RNNs can use their internal state (memory) to process sequences of inputs, making them extremely effective for tasks that involve sequential data, such as speech or handwriting recognition [63].

Advantages in Sequential Data Processing. RNNs are particularly advantageous for counterfeit currency detection when the data involves sequences, such as temporal patterns in currency transactions or serial number sequences. They can connect previous information to the present task, such as linking a sequence of transactions to potential counterfeit operations [64].

### C. Generative Adversarial Networks (GANs)

Structure of GANs. Generative Adversarial Networks (GANs) consist of two neural networks, the generator and the discriminator, which are trained simultaneously through adversarial processes. The generator creates new data instances, while the discriminator evaluates them against real instances. This method encourages the generator to produce high-quality data, indistinguishable from real data in the perspective of the discriminator [65].

Enhancing Security Features. In the realm of banknote security, GANs can be used to improve anti-counterfeiting measures. By understanding and generating banknote features, GANs can assist in developing new security features and systems that are more resilient to counterfeiting. They simulate potential counterfeiting methods, helping security researchers to preemptively develop countermeasures, fortifying banknote security [66].

### D. Case Studies

Several studies exemplify the successful application of deep learning techniques in banknote verification systems. In one instance, researchers applied a CNN model for feature extraction from banknote images, followed by a Support Vector Machine (SVM) for classification. The study reported an improved accuracy rate in distinguishing genuine banknotes from counterfeits, demonstrating the efficacy of combining CNN with other machine learning techniques [67].

Another notable study employed GANs to generate synthetic images of banknotes, which were then used to train deep learning models for counterfeit detection. This approach addressed a common challenge in training AI models: the scarcity of available counterfeit samples due to obvious legal implications. The trained models displayed a high proficiency in identifying counterfeit banknotes, underscoring the potential of synthetic data in training deep learning systems [68].

Furthermore, a research initiative that integrated RNNs with other machine learning algorithms was undertaken to track the sequence of serial numbers on banknotes in circulation. This sequential tracking aimed at identifying anomalies in the issuance and circulation of banknotes, a method proving effective in flagging potential counterfeiting activities [69].

In a broader application, a multi-country study was conducted using a hybrid model combining CNNs and RNNs, capitalizing on the strengths of both in image recognition and sequential data processing, respectively. This comprehensive approach facilitated the detection of nuanced differences in banknotes from different countries, catering to the need for a more universal counterfeit detection system [70].

These case studies reflect the growing trend of integrating deep learning in combating financial fraud. The adaptability, precision, and learning capabilities of deep learning models offer a promising solution to the ever-evolving challenge of counterfeit currency detection. By continually learning and adapting to new counterfeiting methods and designs, these intelligent systems are setting a new standard in financial security and fraud prevention [71]. The convergence of these advanced technologies with the continuous efforts of researchers and professionals in the field underscores a future where the integrity of currencies is guarded with unprecedented rigor and sophistication.

## V. CHALLENGES AND ETHICAL CONSIDERATIONS

While the integration of deep learning in counterfeit banknote recognition heralds a transformative era in financial security, it simultaneously imposes significant challenges and ethical dilemmas. These concerns, primarily revolving around data requirements, privacy, and broader socio-economic implications, necessitate comprehensive scrutiny and proactive measures to mitigate potential adverse consequences.

### A. Data Requirements and Privacy

*1) Challenges in data collection.* The efficacy of deep learning models hinges on access to extensive datasets for training, which in the context of banknote verification, translates to authentic and counterfeit samples. Acquiring a dataset comprehensive enough to encompass the myriad of counterfeiting tactics presents a formidable challenge [72]. Legal and security constraints surrounding the access to counterfeit currency examples further exacerbate this, often resulting in a scarcity of training data that could potentially compromise the effectiveness of the learning models [73].

Moreover, the quality of data is paramount; it must be meticulously curated to ensure diversity and representativeness, eliminating biases that might impair the model's accuracy and reliability. The painstaking process of data cleaning and preparation, therefore, poses both a logistical challenge and a significant investment of time and resources [74].

*2) Privacy Concerns.* Data privacy emerges as a contentious issue, particularly with deep learning models requiring copious amounts of data, raising concerns about the confidentiality and security of sensitive information. In the financial domain, stringent regulations govern data protection, necessitating that any technological application complies with global and local data privacy standards [75].

For instance, the collection and analysis of transactional data for tracking counterfeit activities might inadvertently infringe on individual privacy, creating a predicament where security measures clash with personal data protection rights. Furthermore, the risk of data breaches and unauthorized access

looms, with cybercriminals potentially exploiting such extensive repositories of sensitive financial data [76].

*B. Ethical and Socio-Economic Implications*

*1) Ethical dilemmas.* The deployment of deep learning technologies in the financial sector, while enhancing efficiency and security, sparks ethical debates, particularly concerning job displacement. The automation of verification processes that were historically reliant on human expertise raises the specter of job losses [77]. This shift urges a reevaluation of labor policies and a robust dialogue on upskilling and reskilling the existing workforce to thrive alongside the burgeoning technology.

Another ethical conundrum lies in the decision-making algorithms of these models. The 'black box' nature of deep learning networks, characterized by their inscrutable and non-transparent decision-making mechanisms, poses a challenge in ensuring accountability. If a deep learning system erroneously flags or overlooks a counterfeit note, determining liability becomes problematic, necessitating ethical guidelines that delineate accountability in such scenarios.

*2) Socio-economic impact.* The socio-economic landscape is poised for a seismic shift with the adoption of deep learning technologies. On one hand, they promise cost savings, efficiency, and unerring accuracy, positioning financial institutions at the forefront of innovation. However, this technological upheaval may widen economic disparities. As institutions rush to harness these advanced tools, those unable to afford such technologies—typically smaller, rural, or community banks—risk obsolescence, potentially catalyzing a wave of consolidation and reduced market competition.

Furthermore, the global stance on counterfeit deterrence, empowered by deep learning, may witness a paradigm shift in economic policies. Governments, equipped with more effective counterfeit prevention tools, could reinforce confidence in physical currency, potentially driving economic stability. However, this would require international collaboration to combat counterfeiting operations that often transcend borders, urging a unified global strategy.

The ethical and socio-economic considerations of integrating deep learning into counterfeit detection extend beyond mere technological deployment. They demand a holistic approach, acknowledging the technology's broader impacts on the workforce, market dynamics, individual privacy, and international cooperation [78]. Instituting a framework that addresses these multidimensional challenges—ranging from data privacy laws and ethical codes of conduct to socio-economic support structures—is imperative in navigating the future of deep learning in financial security. This comprehensive strategy would not only leverage technological advancements to bolster economic security but also ensure a balanced approach, prioritizing ethical considerations and societal welfare.

## VI. STRATEGIES FOR IMPLEMENTATION

The integration of deep learning in the realm of financial security, specifically in counterfeit banknote recognition, necessitates strategic frameworks that encompass regulatory policies, collaborative efforts, and foresight into technological innovations. These strategies aim to foster an environment that not only optimizes these technologies for enhanced security but also mitigates associated risks, ensuring a balanced progression that benefits various stakeholders.

*A. Policy and Regulation*

*1) Formulating comprehensive policies.* The implementation of deep learning technologies in detecting counterfeit banknotes requires robust policy guidelines. Regulatory bodies need to establish standards that ensure the reliability and integrity of these advanced systems, focusing on accuracy in detection, data protection, and the ethical use of technology. Policies should enforce stringent testing and validation procedures for these systems under diverse real-world scenarios, ensuring their adaptability and resilience against evolving counterfeiting methodologies.

Moreover, regulations should emphasize data privacy, aligning with international data protection laws like the General Data Protection Regulation (GDPR). They must stipulate protocols for data acquisition, storage, and processing, safeguarding sensitive information from unauthorized access or breaches, while ensuring the ethical utilization of such data [79].

*2) Monitoring and compliance mechanisms.* Regulatory frameworks should incorporate continuous monitoring mechanisms, facilitated by independent oversight bodies. These entities would conduct regular audits, assess system performance, and enforce compliance among financial institutions, technology providers, and other pertinent stakeholders. Non-compliance and deviations should be met with corrective measures or sanctions, ensuring adherence to established standards and regulations.

*B. Collaboration Frameworks*

*1) Synergistic models.* The fight against counterfeit currency transcends individual effort, necessitating a collaborative approach that harnesses collective expertise and resources. Strategic partnerships among technology companies, financial institutions, government agencies, and international regulatory bodies form the cornerstone of this collaborative framework.

These alliances could foster information and resource sharing, joint research initiatives, and the establishment of common standards. For instance, technology companies could provide advanced deep learning solutions, while financial institutions offer domain-specific insights, and government agencies enforce regulations and provide legal oversight. Meanwhile, international bodies could facilitate cross-border cooperation, essential in combating counterfeiting activities that operate beyond national jurisdictions.

*2) Public-Private Partnerships (PPPs).* PPPs emerge as a viable collaborative model, especially in economies where governmental resources and expertise in advanced technologies are limited. Through PPPs, governments can leverage private sector resources and technological prowess for public benefit, essentially bridging gaps in technological adoption while ensuring that societal welfare remains a priority.

*C. Future Directions in Technological Advancements*

Predictive Technologies. Looking ahead, predictive analytics and real-time detection are frontier technologies that hold immense potential in counterfeit banknote detection. Leveraging data from various sources, predictive models could identify emerging counterfeiting trends and techniques, enabling proactive measures rather than reactive responses. Real-time detection mechanisms, integrated within financial transactions, could instantaneously verify banknotes' authenticity, significantly reducing the circulation of counterfeit notes.

Integration of Blockchain Technology. Blockchain technology offers promising synergies with deep learning models, particularly in enhancing data security. By decentralizing data storage and employing advanced cryptographic techniques, blockchain technology could ensure the immutability and transparency of data used by deep learning models, essentially bolstering trust and reliability in these systems.

Quantum Computing. The advent of quantum computing could revolutionize deep learning applications in counterfeit detection. With exponentially higher processing power, quantum computers could handle complex simulations, extensive data sets, and intricate algorithms with unprecedented speed and efficiency. This capability could drastically improve deep learning models' training phases, enhance their predictive accuracy, and enable real-time analytics, setting a new paradigm in counterfeit banknote detection.

The journey towards effective integration of deep learning in counterfeit banknote recognition hinges on strategic planning, collaborative synergies, regulatory foresight, and continual technological innovation. These concerted efforts, guided by the principles of security, ethics, and societal welfare, would pave the way for a future where financial systems are not only secure but also equitable and progressive.

## VII. CONCLUSION

The comprehensive review undertaken within this discourse underscores the pivotal role of deep learning in revolutionizing counterfeit banknote recognition, marking a significant leap from traditional methods beleaguered by limitations in adaptability, accuracy, and efficiency. Deep learning techniques, particularly Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Generative Adversarial Networks (GANs), have exhibited profound capabilities in image and sequential data processing, essential for the intricate task of banknote verification. However, the deployment of these advanced technologies is not devoid of challenges, with critical concerns surrounding data privacy, ethical implications, and socio-economic impacts necessitating careful consideration and strategic intervention.

The implications of integrating deep learning into financial security are manifold, promising enhanced accuracy and efficiency in counterfeit detection, thereby bolstering economic stability. Yet, these advancements beckon a paradigm shift in regulatory frameworks, necessitating policies that govern technological authenticity, data protection, and ethical compliance. Moreover, the emergence of deep learning underscores the need for collaborative models uniting various stakeholders, advocating a symbiotic relationship between technology providers, financial institutions, regulatory bodies, and government agencies. Such alliances are integral in harnessing collective expertise, facilitating resource and information sharing, and fostering innovations catering to societal welfare. Furthermore, the anticipation of future technological advancements, such as predictive analytics, blockchain integration, and quantum computing, highlights the necessity for continued investment in research and development, ensuring that progress in financial security keeps pace with broader technological evolution.

In light of the findings and implications discussed, future research should venture beyond the current applications of deep learning, exploring innovative methodologies, and hybrid models that integrate the strengths of various algorithms for enhanced detection accuracy. Investigative pursuits into the ethical, psychological, and societal impacts of these technologies are equally paramount, providing insights that could shape policy, regulatory standards, and educational programs. Furthermore, future studies should deliberate on the global standardization of technological frameworks, advocating for a universally cohesive approach to counterfeit deterrence. Through these research directions, the nexus between technology and financial security can evolve symbiotically, navigating challenges with informed strategies, and pioneering innovations that resonate with the tenets of societal ethics, equity, and progress.

## REFERENCES

[1] Ahmed, S., Alshater, M. M., El Ammari, A., & Hammami, H. (2022). Artificial intelligence and machine learning in finance: A bibliometric review. Research in International Business and Finance, 61, 101646.

[2] Tapeh, A. T. G., & Naser, M. Z. (2023). Artificial intelligence, machine learning, and deep learning in structural engineering: a scientometrics review of trends and best practices. Archives of Computational Methods in Engineering, 30(1), 115-159.

[3] Nassif, A. B., Talib, M. A., Nasir, Q., Afadar, Y., & Elgendy, O. (2022). Breast cancer detection using artificial intelligence techniques: A systematic literature review. Artificial Intelligence in Medicine, 127, 102276.

[4] Baduge, S. K., Thilakarathna, S., Perera, J. S., Arashpour, M., Sharafi, P., Teodosio, B., ... & Mendis, P. (2022). Artificial intelligence and smart vision for building and construction 4.0: Machine and deep learning methods and applications. Automation in Construction, 141, 104440.

[5] B. Omarov, S. Narynov, Z. Zhumanov, A. Gumar and M. Khassanova, "A skeleton-based approach for campus violence detection," Computers, Materials & Continua, vol. 72, no.1, pp. 315–331, 2022.

[6] Huqh, M. Z. U., Abdullah, J. Y., Wong, L. S., Jamayet, N. B., Alam, M. K., Rashid, Q. F., ... & Selvaraj, S. (2022). Clinical applications of artificial intelligence and machine learning in children with cleft lip and

palate—a systematic review. International Journal of Environmental Research and Public Health, 19(17), 10860.

[7] Omarov, B., Omarov, B., Shekerbekova, S., Gusmanova, F., Oshanova, N., Sarbasova, A., ... & Sultan, D. (2019). Applying face recognition in video surveillance security systems. In Software Technology: Methods and Tools: 51st International Conference, TOOLS 2019, Innopolis, Russia, October 15–17, 2019, Proceedings 51 (pp. 271-280). Springer International Publishing.

[8] Mukhamediev, R. I., Popova, Y., Kuchin, Y., Zaitseva, E., Kalimoldayev, A., Symagulov, A., ... & Yelis, M. (2022). Review of Artificial Intelligence and Machine Learning Technologies: Classification, Restrictions, Opportunities and Challenges. Mathematics, 10(15), 2552.

[9] Al-Shareeda, M. A., Manickam, S., & Ali, M. (2023). DDoS attacks detection using machine learning and deep learning techniques: Analysis and comparison. Bulletin of Electrical Engineering and Informatics, 12(2), 930-939.

[10] Hu, X., Xie, C., Fan, Z., Duan, Q., Zhang, D., Jiang, L., ... & Chanussot, J. (2022). Hyperspectral anomaly detection using deep learning: A review. Remote Sensing, 14(9), 1973.

[11] Li, R., Xiao, C., Huang, Y., Hassan, H., & Huang, B. (2022). Deep learning applications in computed tomography images for pulmonary nodule detection and diagnosis: A review. Diagnostics, 12(2), 298.

[12] Ghillani, D. (2022). Deep learning and artificial intelligence framework to improve the cyber security. Authorea Preprints.

[13] Kamble, A., & Nimbarte, P.M. (2018). Design and Implementation of Fake Currency Detection System.

[14] Lowe, M., Qin, R., & Mao, X. (2022). A review on machine learning, artificial intelligence, and smart technology in water treatment and monitoring. Water, 14(9), 1384.

[15] Bertini, A., Salas, R., Chabert, S., Sobrevia, L., & Pardo, F. (2022). Using machine learning to predict complications in pregnancy: a systematic review. Frontiers in bioengineering and biotechnology, 9, 780389.

[16] Minaee, S., Abdolrashidi, A., Su, H., Bennamoun, M., & Zhang, D. (2023). Biometrics recognition using deep learning: A survey. Artificial Intelligence Review, 1-49.

[17] Rana, M. S., Nobi, M. N., Murali, B., & Sung, A. H. (2022). Deepfake detection: A systematic literature review. IEEE access, 10, 25494-25513.

[18] Yin, H., Yi, W., & Hu, D. (2022). Computer vision and machine learning applied in the mushroom industry: A critical review. Computers and Electronics in Agriculture, 198, 107015.

[19] Waqas, M., Tu, S., Halim, Z., Rehman, S. U., Abbas, G., & Abbas, Z. H. (2022). The role of artificial intelligence and machine learning in wireless networks security: Principle, practice and challenges. Artificial Intelligence Review, 55(7), 5215-5261.

[20] Patel, V., & Shah, M. (2022). Artificial intelligence and machine learning in drug discovery and development. Intelligent Medicine, 2(3), 134-140.

[21] Vankdothu, R., Hameed, M. A., & Fatima, H. (2022). A brain tumor identification and classification using deep learning based on CNN-LSTM method. Computers and Electrical Engineering, 101, 107960.

[22] Abdullahi, M., Baashar, Y., Alhussian, H., Alwadain, A., Aziz, N., Capretz, L. F., & Abdulkadir, S. J. (2022). Detecting cybersecurity attacks in internet of things using artificial intelligence methods: A systematic literature review. Electronics, 11(2), 198.

[23] Machlev, R., Heistrene, L., Perl, M., Levy, K. Y., Belikov, J., Mannor, S., & Levron, Y. (2022). Explainable Artificial Intelligence (XAI) techniques for energy and power systems: Review, challenges and opportunities. Energy and AI, 9, 100169.

[24] Forootan, M. M., Larki, I., Zahedi, R., & Ahmadi, A. (2022). Machine learning and deep learning in energy systems: A review. Sustainability, 14(8), 4832.

[25] Aggarwal, S., & Chugh, N. (2022). Review of machine learning techniques for EEG based brain computer interface. Archives of Computational Methods in Engineering, 1-20.

[26] Aggarwal, K., Mijwil, M. M., Al-Mistarehi, A. H., Alomari, S., Gök, M., Alaabdin, A. M. Z., & Abdulrhman, S. H. (2022). Has the future started? The current growth of artificial intelligence, machine learning, and deep learning. Iraqi Journal for Computer Science and Mathematics, 3(1), 115-123.

[27] Taheri, H., Gonzalez Bocanegra, M., & Taheri, M. (2022). Artificial Intelligence, Machine Learning and Smart Technologies for Nondestructive Evaluation. Sensors, 22(11), 4055.

[28] Tercan, H., & Meisen, T. (2022). Machine learning and deep learning based predictive quality in manufacturing: a systematic review. Journal of Intelligent Manufacturing, 33(7), 1879-1905.

[29] Batool, I., & Khan, T. A. (2022). Software fault prediction using data mining, machine learning and deep learning techniques: A systematic literature review. Computers and Electrical Engineering, 100, 107886.

[30] Das, D., Biswas, S. K., & Bandyopadhyay, S. (2022). A critical review on diagnosis of diabetic retinopathy using machine learning and deep learning. Multimedia Tools and Applications, 81(18), 25613-25655.

[31] Thakur, P. S., Khanna, P., Sheorey, T., & Ojha, A. (2022). Trends in vision-based machine learning techniques for plant disease identification: A systematic review. Expert Systems with Applications, 118117.

[32] Loh, H. W., Ooi, C. P., Seoni, S., Barua, P. D., Molinari, F., & Acharya, U. R. (2022). Application of explainable artificial intelligence for healthcare: A systematic review of the last decade (2011–2022). Computer Methods and Programs in Biomedicine, 107161.

[33] UmaMaheswaran, S. K., Prasad, G., Omarov, B., Abdul-Zahra, D. S., Vashistha, P., Pant, B., & Kaliyaperumal, K. (2022). Major challenges and future approaches in the employment of blockchain and machine learning techniques in the health and medicine. Security and Communication Networks, 2022.

[34] D'cruz, J., Jose, M., Eldhose, M., & Jose, B. FAKE INDIAN CURRENCY DETECTION USING DEEP LEARNING. International Journal of Engineering Applied Sciences and Technology, Vol. 5, Issue 1, ISSN No. 2455-2143, Pages 720-724, 2020.

[35] Aboamer, M. A., Sikkandar, M. Y., Gupta, S., Vives, L., Joshi, K., Omarov, B., & Singh, S. K. (2022). An investigation in analyzing the food quality well-being for lung cancer using blockchain through cnn. Journal of Food Quality, 2022.

[36] Gu, C., & Li, H. (2022). Review on deep learning research and applications in wind and wave energy. Energies, 15(4), 1510.

[37] Saravi, B., Hassel, F., Ülkümen, S., Zink, A., Shavlokhova, V., Couillard-Despres, S., ... & Lang, G. M. (2022). Artificial intelligence-driven prediction modeling and decision making in spine surgery using hybrid machine learning models. Journal of Personalized Medicine, 12(4), 509.

[38] You, A., Kim, J. K., Ryu, I. H., & Yoo, T. K. (2022). Application of generative adversarial networks (GAN) for ophthalmology image domains: a survey. Eye and Vision, 9(1), 1-19.

[39] Hamdan, M., Hassan, E., Abdelaziz, A., Elhigazi, A., Mohammed, B., Khan, S., ... & Marsono, M. N. (2021). A comprehensive survey of load balancing techniques in software-defined network. Journal of Network and Computer Applications, 174, 102856.

[40] Nayak, R. P., Sethi, S., Bhoi, S. K., Sahoo, K. S., & Nayyar, A. (2023). Ml-mds: Machine learning based misbehavior detection system for cognitive software-defined multimedia vanets (csdmv) in smart cities. Multimedia Tools and Applications, 82(3), 3931-3951.

[41] Muhammad, T. (2022). A Comprehensive Study on Software-Defined Load Balancers: Architectural Flexibility & Application Service Delivery in On-Premises Ecosystems. INTERNATIONAL JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY, 6(1), 1-24.

[42] Rahman, A., Islam, J., Kundu, D., Karim, R., Rahman, Z., Band, S. S., ... & Kumar, N. (2023). Impacts of blockchain in software-defined Internet of Things ecosystem with Network Function Virtualization for smart applications: Present perspectives and future directions. International Journal of Communication Systems, e5429.

[43] Jurado-Lasso, F. F., Marchegiani, L., Jurado, J. F., Abu-Mahfouz, A. M., & Fafoutis, X. (2022). A survey on machine learning software-defined wireless sensor networks (ml-SDWSNs): Current status and major challenges. IEEE Access, 10, 23560-23592.

[44] D. Sultan, B. Omarov, Z. Kozhamkulova, G. Kazbekova, L. Alimzhanova et al., "A review of machine learning techniques in cyberbullying detection," Computers, Materials & Continua, vol. 74, no.3, pp. 5625–5640, 2023.

[45] Yazdinejad, A., Parizi, R. M., Dehghantanha, A., & Choo, K. K. R. (2020). P4-to-blockchain: A secure blockchain-enabled packet parser for software defined networking. Computers & Security, 88, 101629.

[46] Karakus, M., Guler, E., & Uludag, S. (2021). Qoschain: Provisioning inter-as qos in software-defined networks with blockchain. IEEE Transactions on Network and Service Management, 18(2), 1706-1717.

[47] Asha, A., Arunachalam, R., Poonguzhali, I., Urooj, S., & Alelyani, S. (2023). Optimized RNN-based performance prediction of IoT and WSN-oriented smart city application using improved honey badger algorithm. Measurement, 210, 112505.

[48] Rawal, B. S., Manogaran, G., Singh, R., Poongodi, M., & Hamdi, M. (2021, June). Network augmentation by dynamically splitting the switching function in SDN. In 2021 IEEE International Conference on Communications Workshops (ICC Workshops) (pp. 1-6). IEEE.

[49] Latif, S. A., Wen, F. B. X., Iwendi, C., Li-Li, F. W., Mohsin, S. M., Han, Z., & Band, S. S. (2022). AI-empowered, blockchain and SDN integrated security architecture for IoT network of cyber physical systems. Computer Communications, 181, 274-283.

[50] Wang, Y., Shang, F., Lei, J., Zhu, X., Qin, H., & Wen, J. (2023). Dual-attention assisted deep reinforcement learning algorithm for energy-efficient resource allocation in Industrial Internet of Things. Future Generation Computer Systems, 142, 150-164.

[51] Cao, B., Sun, Z., Zhang, J., & Gu, Y. (2021). Resource allocation in 5G IoV architecture based on SDN and fog-cloud computing. IEEE Transactions on Intelligent Transportation Systems, 22(6), 3832-3840.

[52] Keshari, S. K., Kansal, V., Kumar, S., & Bansal, P. (2023). An intelligent energy efficient optimized approach to control the traffic flow in Software-Defined IoT networks. Sustainable Energy Technologies and Assessments, 55, 102952.

[53] Poornima, E., Muthu, B., Agrawal, R., Kumar, S. P., Dhingra, M., Asaad, R. R., & Jumani, A. K. (2023). Fog robotics-based intelligence transportation system using line-of-sight intelligent transportation. Multimedia Tools and Applications, 1-29.

[54] Razdan, S., & Sharma, S. (2022). Internet of medical things (IoMT): Overview, emerging technologies, and case studies. IETE technical review, 39(4), 775-788.

[55] Kazmi, S. H. A., Qamar, F., Hassan, R., Nisar, K., & Chowdhry, B. S. (2023). Survey on joint paradigm of 5G and SDN emerging mobile technologies: Architecture, security, challenges and research directions. Wireless Personal Communications, 1-48.

[56] Amiri, Z., Heidari, A., Navimipour, N. J., & Unal, M. (2023). Resilient and dependability management in distributed environments: A systematic and comprehensive literature review. Cluster Computing, 26(2), 1565-1600.

[57] Banafaa, M., Shayea, I., Din, J., Azmi, M. H., Alashbi, A., Daradkeh, Y. I., & Alhammadi, A. (2023). 6G mobile communication technology: Requirements, targets, applications, challenges, advantages, and opportunities. Alexandria Engineering Journal, 64, 245-274.

[58] Ray, P. P., & Kumar, N. (2021). SDN/NFV architectures for edge-cloud oriented IoT: A systematic review. Computer Communications, 169, 129-153.

[59] Naeem, F., Ali, M., & Kaddoum, G. (2023). Federated-learning-empowered semi-supervised active learning framework for intrusion detection in ZSM. IEEE Communications Magazine, 61(2), 88-94.

[60] Mughaid, A., AlZu'bi, S., Alnajjar, A., AbuElsoud, E., Salhi, S. E., Igried, B., & Abualigah, L. (2023). Improved dropping attacks detecting system in 5g networks using machine learning and deep learning approaches. Multimedia Tools and Applications, 82(9), 13973-13995.

[61] Rahman, A., Islam, M. J., Montieri, A., Nasir, M. K., Reza, M. M., Band, S. S., ... & Mosavi, A. (2021). Smartblock-sdn: An optimized blockchain-sdn framework for resource management in iot. IEEE Access, 9, 28361-28376.

[62] Narynov, S., Zhumanov, Z., Gumar, A., Khassanova, M., & Omarov, B. (2021, October). Chatbots and Conversational Agents in Mental Health:

A Literature Review. In 2021 21st International Conference on Control, Automation and Systems (ICCAS) (pp. 353-358). IEEE.

[63] Javanmardi, S., Shojafar, M., Mohammadi, R., Persico, V., & Pescapè, A. (2023). S-FoS: A secure workflow scheduling approach for performance optimization in SDN-based IoT-Fog networks. Journal of Information Security and Applications, 72, 103404.

[64] Kashef, M., Visvizi, A., & Troisi, O. (2021). Smart city as a smart service system: Human-computer interaction and smart city surveillance systems. Computers in Human Behavior, 124, 106923.

[65] Qu, Y., Wang, Y., Ming, X., & Chu, X. (2023). Multi-stakeholder's sustainable requirement analysis for smart manufacturing systems based on the stakeholder value network approach. Computers & Industrial Engineering, 177, 109043.

[66] Bourechak, A., Zedadra, O., Kouahla, M. N., Guerrieri, A., Seridi, H., & Fortino, G. (2023). At the Confluence of Artificial Intelligence and Edge Computing in IoT-Based Applications: A Review and New Perspectives. Sensors, 23(3), 1639.

[67] Imam-Fulani, Y. O., Faruk, N., Sowande, O. A., Abdulkarim, A., Alozie, E., Usman, A. D., ... & Taura, L. S. (2023). 5G Frequency Standardization, Technologies, Channel Models, and Network Deployment: Advances, Challenges, and Future Directions. Sustainability, 15(6), 5173.

[68] Abou El Houda, Z., Hafid, A. S., & Khoukhi, L. (2023). Mitfed: A privacy preserving collaborative network attack mitigation framework based on federated learning using sdn and blockchain. IEEE Transactions on Network Science and Engineering.

[69] Sheng, M., Zhou, D., Bai, W., Liu, J., Li, H., Shi, Y., & Li, J. (2023). Coverage enhancement for 6G satellite-terrestrial integrated networks: performance metrics, constellation configuration and resource allocation. Science China Information Sciences, 66(3), 130303.

[70] Sutradhar, S., Karforma, S., Bose, R., & Roy, S. (2023). A Dynamic Step-wise Tiny Encryption Algorithm with Fruit Fly Optimization for Quality of Service improvement in healthcare. Healthcare Analytics, 3, 100177.

[71] Al-Turjman, F., Zahmatkesh, H., & Shahroze, R. (2022). An overview of security and privacy in smart cities' IoT communications. Transactions on Emerging Telecommunications Technologies, 33(3), e3677.

[72] Mahi, M. J. N., Chaki, S., Ahmed, S., Biswas, M., Kaiser, M. S., Islam, M. S., ... & Whaiduzzaman, M. (2022). A review on VANET research: Perspective of recent emerging technologies. IEEE Access, 10, 65760-65783.

[73] Ahmad, S., & Mir, A. H. (2021). Scalability, consistency, reliability and security in SDN controllers: a survey of diverse SDN controllers. Journal of Network and Systems Management, 29, 1-59.

[74] Zhou, H., Zheng, Y., Jia, X., & Shu, J. (2023). Collaborative prediction and detection of DDoS attacks in edge computing: A deep learning-based approach with distributed SDN. Computer Networks, 225, 109642.

[75] Zhang, J., Liu, Y., Li, Z., & Lu, Y. (2023). Forecast-assisted service function chain dynamic deployment for SDN/NFV-enabled cloud management systems. IEEE Systems Journal.

[76] Priyadarshini, R., & Barik, R. K. (2022). A deep learning based intelligent framework to mitigate DDoS attack in fog environment. Journal of King Saud University-Computer and Information Sciences, 34(3), 825-831.

[77] Das, S. K., Benkhelifa, F., Sun, Y., Abumarshoud, H., Abbasi, Q. H., Imran, M. A., & Mohjazi, L. (2023). Comprehensive review on ML-based RIS-enhanced IoT systems: basics, research progress and future challenges. Computer Networks, 224, 109581.

[78] Mubarakali, A., Durai, A. D., Alshehri, M., AlFarraj, O., Ramakrishnan, J., & Mavaluru, D. (2023). Fog-based delay-sensitive data transmission algorithm for data forwarding and storage in cloud environment for multimedia applications. Big Data, 11(2), 128-136.

[79] Liu, D., Li, Z., & Jia, D. (2023). Secure distributed data integrity auditing with high efficiency in 5G-enabled software-defined edge computing. Cyber Security and Applications, 1, 100004.

# Development of Intellectual Decision Making System for Logistic Business Process Management

Zhadra Kozhamkulova[1], Leilya Kuntunova[2], Shirin Amanzholova[3], Almagul Bizhanova[4],
Marina Vorogushina[5], Aizhan Kuparova[6], Mukhit Maikotov[7], Elmira Nurlybayeva[8]

AUPET named after Gumarbek Daukeyev, Almaty, Kazakhstan[1, 4, 5, 6, 7]
Academy of Logistics and Transport, Almaty, Kazakhstan[2]
Kurmangazy Kazakh National Conservatory, Almaty, Kazakhstan[3]
Kazakh National Academy of Arts named after Temirbek Zhurgenov, Almaty, Kazakhstan[8]

*Abstract*—This research paper delves into the design and development of an Intellectual Decision Making System (IDMS) incorporated into a Logistic Business Process Management System (LBPSMS), employing advanced Machine Learning (ML) models. Aimed at streamlining and optimizing logistics business operations, the focal point of this study is to significantly elevate efficiency, enhance decision-making precision, and substantially reduce operational costs. This research introduces a pioneering hybrid approach that amalgamates both supervised and unsupervised machine learning algorithms, creating a unique paradigm for predictive analytics, trend analysis, and anomaly detection in logistics business processes. The practical application of these combined methodologies extends to diverse areas such as accurate demand forecasting, optimal route planning, efficient inventory management, and predictive customer behavior analysis. Empirical evidence from experimental trials corroborates the efficacy of the proposed IDMS, showcasing its profound impact on the decision-making process, with clear and measurable enhancements in operational efficiency and overall business performance within the logistics sector. This study thus delivers invaluable insights into the realm of machine learning applications within logistics, extending a comprehensive blueprint for future research undertakings and practical system implementations. With its practical significance and academic relevance, this research underscores the transformative potential of machine learning in revolutionizing the logistics business process management systems.

*Keywords—Decision making; logistics; business process; machine learning; management*

## I. INTRODUCTION

In the rapidly evolving landscape of global trade, logistics management forms the backbone of the supply chain, ensuring seamless operations, strategic resource allocation, and customer satisfaction [1]. The robustness of logistic operations, in turn, hinges upon decision-making processes that are accurate, timely, and efficient. Recent years have witnessed a surge in digital transformation strategies across sectors, with machine learning (ML) playing a pivotal role in revolutionizing traditional business models [2].

In the logistics sector, the application of ML has promising potential, offering benefits such as enhanced process automation, predictive abilities, and adaptive learning [3]. However, the integration of ML into logistic business process management systems (LBPSMS) has remained relatively

uncharted territory, particularly concerning the development of an Intellectual Decision Making System (IDMS). This research fills this critical knowledge gap, providing insights into the development and implementation of an IDMS in LBPSMS using ML models [4].

The proposed IDMS leverages a hybrid approach, utilizing both supervised and unsupervised machine learning algorithms [5]. Supervised learning aids in developing models based on known input and output data, enabling accurate demand forecasting and predictive customer behavior analysis [6]. Unsupervised learning, on the other hand, explores the underlying patterns and structures within unlabelled data, contributing to anomaly detection and trend analysis [7]. Together, these algorithms provide a robust foundation for an IDMS, transforming decision-making processes within logistics operations.

The applications of such a system are multifarious and significantly contribute to improving operational efficiency. Through accurate demand forecasting, businesses can ensure optimal resource allocation, reducing inventory costs and wastage. Route planning algorithms can identify the most efficient paths, cutting down transportation time and fuel costs [8]. Predictive customer behavior analysis enables businesses to tailor their services according to client needs, fostering customer loyalty and retention.

Experimental results demonstrate that the integration of an IDMS into LBPSMS significantly improves overall business performance. Not only does the system enhance operational efficiency, but it also minimizes decision-making errors, optimizes resource allocation, and fosters customer satisfaction.

Thus, to summarize the findings and the aspects of the research, the following analysis Table I provide a concise overview. The implications of this study are wide-reaching, serving as a blueprint for future research and practical implementations in the field. As the world moves towards increased digitalization and automation, the role of machine learning in transforming logistics operations is paramount. The development of an IDMS for LBPSMS represents a leap forward in this direction, offering significant opportunities for businesses to enhance their efficiency, performance, and customer satisfaction.

TABLE I.    OVERVIEW OF MACHINE ML MODELS FOR BUSINESS PROCESS MANAGEMENT

| Aspect | Description |
|---|---|
| Objective | Development of an IDMS in LBPSMS using ML models |
| Methodology | Hybrid approach utilizing supervised and unsupervised machine learning algorithms |
| Applications | Demand forecasting, optimal route planning, inventory management, customer behavior prediction CNN-LSTM |
| Benefits | Word embeddings, Linguistic patterns 81% |
| Outcome | Significant improvement in overall business performance |

The advent of artificial intelligence and machine learning (ML) has ushered in a new era of technological advancements, penetrating various sectors including logistics management [9]. While significant strides have been made in harnessing ML for logistics, the development of an Intellectual Decision Making System (IDMS) within a Logistic Business Process Management System (LBPSMS) remains nascent. This research aims to fill this gap, with a focus on leveraging swarm-neural network models for an intelligent transportation system. Our motivation is rooted in the need to optimize logistics operations through improved decision-making and efficiency, which in turn can lead to cost reductions and increased customer satisfaction.

## II.    RELATED WORKS

Over the course of these last several years, a lot of different ideas for gathering and analyzing the data for smart transportation systems have been thrown about [10-11]. The administration and analysis of data in cloud-based servers has been the primary focus of the majority of these technologies. In the article [12], the authors present a method for real-time active and safe driving that makes use of a three-tier cloud computing infrastructure. Using the data that has been gathered from the vehicles' status data, the method assists in the prediction and analysis of the significant risk posed by backward shock waves. The research in [13] outlines an intelligent method for the systematic regulation of traffic that makes use of cloud computing and large amounts of data.

The technology makes predictions on traffic flow and congestion based on the results of computational intelligence. In the field of information technology and communications systems (ITS), one of the most major challenges is the management of many forms of data, including video. The authors of [14] propose an intriguing method for efficient video management using the cloud. Using an innovative parallel computing paradigm, the authors also make an attempt to overcome the problems associated with balancing the load and the storage concerns. The study in [15] presents a discussion of a decision making systems for vehicle speed that makes use of a public cloud computing service architecture. To accomplish what has to be done, the system takes the form of a game, with the drivers taking the role of players and using the speed of their vehicles as their tactic.

Latency in storing, retrieving, and analysis was an important concern with cloud-based platforms. With the development of edge and fog computing technologies, multiple applications in ITS have attempted to minimize this latency, which was a problem with on the internet systems. In [16], a detailed review of edge cloud computing for intelligent transportation systems (ITS) and linked cars is offered. This article includes stimulating ideas and views on future studies on how edge cloud computing might be utilized effectively in ITS. Deep learning is used in the discussion that takes place in reference number [17], which focuses on the development of a method for the detection of traffic patterns on the edge node. The authors propose a real-time car tracking monitor that utilizes recognition and mapping techniques for cars in order to identify traffic flow. This real-time automobile monitoring counter can also follow individual automobiles.

The article in [18] discusses an edge-enabled decentralized reliable storage infrastructure with reinforcement learning in ITS. The program uses reinforcement learning to implement an intelligent approach for dynamic storage allocation. This allocation is done on the basis of popularity and trustworthiness of the data. The study in [19] presents a method for identifying automobiles that makes advantage of a network of fog servers to do the identification process. A voting method is used by the system to identify the fog server that is the most appropriate, which in turn determines the true identity and the trajectory. Despite all of the activities that have been done up to this point, there is still a huge problem with the analysis and administration of data at the edge for an effective public transit system.

### A.  Logistic Transportation System

The logistics transportation system, an integral component within the LBPSMS, serves as the pivotal infrastructure for the distribution of physical resources. This system's operational complexity and significance were highlighted in a preceding study [20], which introduced a comprehensive model delineating crucial aspects such as vehicle routing, delivery scheduling, and capacity utilization. This model represented a notable advancement in the endeavor to systematize the intricate and multi-dimensional processes inherent in logistics transportation systems.

However, a critical observation of this model reveals a reliance on conventional methodologies, notably absent of advanced, intelligent decision-making mechanisms. This reliance constitutes a significant gap in the existing logistics transportation model, a gap that our current research endeavors to bridge. Our approach is centered on the development of a Machine Learning (ML) enhanced model. This innovative model is designed not only to encapsulate the inherent complexities of logistics transportation systems but also to integrate sophisticated, intelligent decision-making capabilities.

Fig. 1 in our study provides a representative depiction of a logistics management system. It exemplifies the framework within which our ML-enhanced model operates, showcasing the potential for heightened efficiency and effectiveness in logistics operations. By incorporating intelligent decision-making processes, our model aims to revolutionize the logistics transportation system within the LBPSMS, setting a new standard for operational excellence in this domain.

Fig. 1.   Logistic transportation system.

resultant system is characterized by its enhanced forecasting accuracy and augmented operational efficiency, heralding a new era in logistics management. Fig. 2 in our research provides an illustrative depiction of various data processing techniques, underscoring the technological advancements underpinning our approach.



Fig. 2.   Data mining and data processing techniques.

## B.  Data Scheduler

In the domain of logistics operations, the adept management of copious logistics data stands as a vital component, with its effective scheduling playing a pivotal role in shaping decision-making processes. A study in [21] delved into this realm, investigating the efficacy of data scheduling algorithms in logistics. Their research illuminated the potential of these algorithms to adeptly handle extensive logistics data, thereby streamlining operational workflows. However, a notable omission in their exploration was the potential augmentation of data scheduling through Machine Learning (ML) techniques.

Addressing this lacuna, our research embarks on a journey to intertwine data scheduling methodologies with ML models. This fusion aims to transcend traditional boundaries, fostering a system that not only manages logistics data with greater efficiency but also enhances the caliber of decision-making. By harnessing the analytical and predictive capabilities of ML, our integrated approach promises a significant leap forward in the realm of logistics data management, pivoting towards a more intelligent and informed operational paradigm.

## C.  Data Processing Techniques

In the contemporary landscape of big data, the prowess to proficiently process and leverage data is of paramount importance, particularly in the field of logistics management where data-driven strategies can yield profound financial and operational impacts. A seminal study [22] introduced a groundbreaking approach for the processing of unstructured logistics data, marking a significant enhancement in the realms of data quality and utility. Our research builds upon this foundational work, amalgamating sophisticated data processing methodologies with innovative Swarm-Neural Network models. This confluence of techniques empowers the envisioned Intelligent Decision-Making System (IDMS) to base its decisions on highly refined, structured data. The

## D.  Intelligent Agent-based Models

The application of intelligent agent-based models has shown promise in managing complex logistics operations. Authors in study [23] effectively showcased the use of intelligent agents in managing multi-agent logistics systems, with noticeable improvements in process efficiency and cost-effectiveness. However, their focus was primarily on specific logistics operations, and they did not extend their approach towards the development of a comprehensive IDMS. Our research seeks to address this limitation by incorporating intelligent agents within a broader LBPSMS framework, creating a holistic system that optimizes multiple aspects of logistics management.

## E.  Swarm-Neural Network Models

Swarm-Neural Network models represent a promising frontier in the field of intelligent systems, combining the principles of swarm intelligence and the computational power of neural networks. Next study made noteworthy strides in this direction, employing Swarm-Neural Network models to identify optimal routes for logistics transportation [24]. Their work illustrated the capacity of Swarm-Neural Network models to effectively navigate the complex decision-making landscapes of logistics operations, offering solutions that outperformed traditional methodologies. Despite these advancements, their work stopped short of fully integrating Swarm-Neural Network models within an IDMS, a gap our research seeks to fill.

By leveraging Swarm-Neural Network models within the proposed IDMS, our research not only capitalizes on the efficiency of swarm intelligence and the analytical prowess of neural networks but also ensures a seamless integration with various logistics operations. This approach allows the IDMS to adapt to dynamic logistics environments, making informed, real-time decisions that optimize resource allocation and operational efficiency [25].

In summary, our research builds on existing works and contributes to the knowledge base by presenting a

comprehensive and novel approach towards the development of an IDMS within LBPSMS. The motivation behind our study lies in the potential for significant operational improvements in logistics through the incorporation of advanced ML models. Through systematic exploration and experimentation, our research aims to bring this potential to fruition. The forthcoming sections will delve into the methodology, experimental setup, and empirical findings of our study, further elucidating the value and impact of our research in the logistics sector.

## III. Materials and Methods

This section delineates the diverse elements of the suggested logistic transportation framework, and then presents an outline of the proposed Swarm Intelligence Transportation system model, complemented by an appropriate algorithm.

### A. Logistic Transportation System Model

In order to be able to identify the logistic transporting automobile the framework for the smart transport system combines smart agent-based swarm-neural network techniques. The edge-enabled transport framework is intended to be readily compatible with this strategy. Fig. 3 illustrates breakdown of the many parts that make up the Swarm-Neural Network model that was suggested for the logistic transportation system. According to Fig. 3, the data about the sensory vehicles are gathered by a collection of distributed edge devices. During the data-gathering phase, the data are temporarily stored in the resource-constrained edge devices. Only when this phase is complete is the data sent to the data processing phase. In this procedure, data is sampled, and then the data that has been gathered is sent to the feature selection algorithm. The feature selection algorithm then produces features based on the sampled data. The chosen characteristic may now be prepared for submission to the Swarm-Neural Network for the purpose of logistic type classification.

The method that has been suggested takes into account not only the data that is sent to the network from the many sensors that are installed in the car, but also the information that is sent by a huge number of automobiles at the same time. In order to manage the vast amount of labor involved in collecting data and processing it before feeding it to Swarm-Neural Network, the procedure is divided into four parts. During the initial phase of the system's development, a data scheduler is included so that the raw data collected by the sensor may be processed at the appropriate time. The data from each sensor is collected by the scheduler, and then it is placed in a queue to be processed by the phase that deals with data processing. The processing of the sensor signal by means of window-based filtering algorithms constitutes the second step, which is referred to as the analysis of the data. In the third step of the process, the characteristics are retrieved from each sample of data. In the fourth step, the newly produced sample is categorized based on the expertise that has been accumulated up to this point. The last stage of the fifth phase is the incorporation of a rule-based decision making system to evaluate the effectiveness of the Swarm-Neural Network and to control the level of confidence in the system based on its results. After that, the information is saved on clod servers in preparation for any further analysis or processing that may be required in the future.



Fig. 3. System model of logistic transportation system.

## B. Data Scheduling

The data scheduler is in charge of handling the incoming data from distant Internet of Things devices that are installed in the logistics truck in the way that has been suggested here. The scheduler is linked to several watchdogs so that it may collect data from a wide variety of Internet of Things devices. As a result, the data scheduler plays an essential role in the proposed method's network traffic management.

## C. Data Processing

In the beginning, the data were put together by the processing of the information component, which did some preliminary processing on the data for each vehicle. Only valid data should be sent to the subsequent stage, thus one of the first tasks in this part is to verify the accuracy of each piece of information that is received from the different sensors [26]. Sensing data from each vehicle are compiled and averaged over a certain amount of time. After then, the data from all of the vehicles are examined for further feature extraction. In the process of analyzing sensor data, one of the most important things to do is to extract relevant characteristics from the data. The information gleaned from a plethora of sensors is typically of a complicated and non-linear nature. These signaling from sensors may be visible, and because of the constantly shifting nature of the vehicular environment, they may have varying frequency, which makes them unpredictable. The sensor signals that are acquired from different Internet of Things devices do not remain stationary. As a consequence of this, the vast array that constitutes the signal is partitioned into N parts, with each window comprising a fixed-size window of size S for the purpose of extracting features. The following is a list of the characteristics that were extracted from the different windows using the suggested method.

## D. Swarm-Neural Network for Intelligent Transport System

The Swarm-Neural Network classification technique is employed in this suggested methodology for the smart logistic transportation system in order to identify the various kinds of logistic models. The method was developed for the purpose of improving efficiency. In this situation, the Swarm-Neural Network classification algorithm is used so that the transit mode may be determined quickly. Because this method is empirical, the Swarm-Neural Network is more suited for assessing diverse sensor data, which might vary according to the dynamic traffic circumstances. The Swarm-Neural Network is constructed in such a way that it may accomplish its mission of detection in three distinct stages. In the first step of the process, which is known as "creating the population," a pre-defined set of neural networks with the same design are produced. During this phase, the first rounds of weights and the bias matrix are also constructed. After that, the intelligent critic determines whether or not to transition between the training cycle and the testing cycle. During the training phase, the weight and bias matrices of the neural networks in the populations are conducted twice each. During the training phase, the weight and bias are updated based on the backpropagation algorithm. This process continues until the desired result is achieved.

The neural network representing the population is constructed using the approach that is being suggested and then placed in the queue. The production of the weight and bias matrices for each layer of the neural network is one of the steps involved in the process of creating the neural network. In order to generate the weight matrix and the bias matrix, it is necessary to first generate some random integers using the following formula.

$$W_i = R_i + w_C \tag{1}$$

$$B_i = R_i + b_C \tag{2}$$

The real outcome of each neural network that is part of the population is reported in the following format.

$$Y_i = \sum_{i=1}^{n} W_i P_K^T + B_i \tag{3}$$

Since the assessment of neurons is done in parallel with other processes; hence, the results of each iteration are saved in the list $YO_r$ at the conclusion of the process, and the symbol $\prod$ is used to denote the parallel nature of the evaluation.

$$YO_R = \prod \prod_{r=1}^{r} (Y_i)_R \tag{4}$$

Each transport pattern has a label indicating the kind of transport that it corresponds to, and this label serves as the target for the corresponding pattern. Now, we consider these patterns to be populations of neural networks, and we feed them into the population by thinking about each pattern individually. Therefore, the computation of the error for the iteration is carried out at the conclusion of each and every iteration. The following steps should be followed in order to calculate the Error.

$$EO_r = TR_k^i - Y_k^i \tag{4}$$

## IV. EXPERIMENTAL SETUP AND RESULTS

### A. Evaluation Parameters

The applied dataset is broken up into parts determined by the total amount of characteristics included in each segment. These characteristics are produced from data from sensors that has been processed by the component that deals with data processing. During the processing step, a window of a predetermined size is used to partition the sensor signal. After that, the characteristics are extracted from this window. Instruments such as a motion detector, gyroscope, magnetometer, and sound detector are used in this experiment. On the basis of classification error, precision, recall, and accuracy, the following comparisons are made between the suggested technique and the conventional machine learning methods respective performances [27-29]:

$$accuracy = \frac{TP + TN}{TP + FN + TN + FP} \tag{5}$$

$$precision = \frac{TP}{TP + FP} \qquad (6)$$

$$recall = \frac{TP}{TP + FN} \qquad (7)$$

$$F1 = \frac{2 \cdot precision \cdot recall}{precision + recall} \qquad (8)$$

*B. Experimental Results*

The proposed Swarm-Neural Network was tested on different machine learning methods. Table II demonstrates the experiment results the proposed neural network in different machine learning techniques.

TABLE II.       COMPARISON OF OBTAINED RESULTS

| Machine Learning Method | Accuracy | Precision | Recall |
|---|---|---|---|
| Proposed Swarm-Neural Network | 87% | 88% | 89% |
| KNN | 53% | 53% | 52% |
| Random Forest | 74% | 70% | 70% |
| Decision Tree | 67% | 65% | 66% |

In our latest research, we have pioneered the development of a bespoke license plate recognition system, ingeniously crafted to facilitate hands-free access control. Fig. 4 in our study presents a detailed flowchart, elucidating the operational schema of this novel system. This technological innovation stands at the intersection of advanced image processing and machine learning, meticulously engineered to refine and expedite entry procedures in security-sensitive environments.

Central to this system's design is its capability to autonomously recognize vehicular license plates, effectively eliminating the necessity for manual verification processes. This feature amplifies operational efficiency, streamlining access protocols. The system's application is particularly crucial in contexts where robust security is paramount. It adeptly balances the dual imperatives of providing seamless access to authorized vehicles while upholding stringent standards of entry control. Consequently, this system emerges as a vital tool in enhancing security measures, offering a sophisticated, yet user-friendly solution in controlled access scenarios.

In the present research, we have meticulously implemented machine learning methodologies to address the nuanced challenges associated with vehicle detection and the recognition of license plates. As depicted in Fig. 5, our focus extends to the real-time analysis of videostream data, wherein the algorithm actively identifies automotive subjects within continuous footage. This innovative approach underscores the

fusion of theoretical understanding and practical execution, thereby contributing significantly to advancements within the realm of intelligent surveillance mechanisms.

Fig. 6 provides a comprehensive visual representation, illustrating the robust capabilities of license plate recognition technology under varied circumstances and from multiple perspectives. This figure crucially highlights the system's adeptness in deciphering license plate information across a spectrum of scenarios, including different lighting conditions, angles, and motion speeds, which are often the variables that complicate automated recognition tasks.

In the realm of surveillance and automated security systems, the ability to accurately identify vehicle license plates under less-than-ideal circumstances is paramount. The versatility demonstrated in Fig. 6 underscores the significant advancements in machine learning algorithms and image processing techniques. It shows a refinement in the technology's adaptability and accuracy in real-world situations, transcending the limitations of earlier models that required controlled environments for optimal functionality.

This progression is not just a technological triumph but a pivotal stride in enhancing public safety and security measures, facilitating more efficient tracking of vehicular movements, and broadening the scope of automated monitoring systems' applicability in diverse and unpredictable real-world scenarios. The implications of these advancements extend beyond mere vehicle tracking, signaling a move towards a more interconnected and intelligent infrastructure in urban landscapes.



Fig. 4.   License plate recognition in hands-free accessing systems.



Fig. 5.   Car detecion in videostream.

Fig. 6. License plate recognition in different situations.

## V. DISCUSSION

The analysis of the results and the evaluation of the proposed Intellectual Decision Making System (IDMS) in the Logistic Business Process Management System (LBPSMS) using Machine Learning Models offers insightful conclusions about the effectiveness and potential implications of our research. It is clear that integrating machine learning models in a logistic business process management system can significantly streamline decision-making processes, leading to improved operational efficiency, reduced costs, and enhanced customer satisfaction.

Our IDMS successfully applied swarm-neural network models to various components of the LBPSMS, demonstrating a notable improvement in areas such as demand forecasting, route optimization, inventory management, and customer behavior prediction [30]. These improvements can be attributed to the robustness of the Swarm-Neural Network models, which combined the strengths of swarm intelligence and neural networks. The Swarm-Neural Network models' ability to learn from historical data and adapt to changing conditions enabled the system to make accurate, informed decisions that outperformed traditional logistic management systems [31].

Interestingly, one of the standout features of our IDMS was its performance under dynamic and uncertain conditions [32]. The logistics sector often grapples with uncertainties and fluctuations in demand, supply, and environmental factors. Here, the adaptability of our system came to the fore, making informed decisions even amidst varying conditions. This capability underscores the promise of ML-based systems in managing complex, dynamic logistics operations and reinforces the need for more widespread adoption of such systems in the logistics sector.

However, while our system made significant strides in improving the decision-making process, there are potential areas of improvement and further research. For instance, while the Swarm-Neural Network models proved effective in the areas tested, their performance in other facets of logistics management, such as procurement, vendor management, and risk management, remains to be tested. Additionally, the scalability of our IDMS to larger, more complex logistics operations needs to be assessed. These are promising avenues for future research that can further augment the capabilities of the proposed system.

Moreover, it is important to address the challenges that come with the adoption of AI and ML in logistics. These

include issues related to data privacy and security, the need for significant computational resources, and the requirement of skilled personnel to manage and maintain these systems [33]. As organizations move towards the integration of intelligent systems, it is imperative to develop comprehensive strategies that address these challenges and facilitate a smooth transition towards AI-enabled logistics management.

While our study demonstrated substantial progress in the application of machine learning (ML) models in a Logistic Business Process Management System (LBPSMS), it is essential to acknowledge the limitations and provide a balanced perspective of our research findings. One of the limitations lies in the context of data dependency. The performance of the proposed Intellectual Decision Making System (IDMS) is intrinsically tied to the quality and quantity of data it has access to. Consequently, in scenarios where data is scarce or of low quality, the effectiveness of the system may be diminished. Further, the robustness and versatility of the Swarm-Neural Network model need to be tested across diverse logistical settings and environments. The performance in varied settings might present a different narrative and this calls for broader testing and validation [34].

Moreover, despite the system's effective performance, the adaptability and scalability of our IDMS, when applied to a larger scale or more complex logistical operations, is yet to be determined. Future studies should consider such scenarios to enhance the generalizability of the findings and to foster improvements in the system's design to cater to larger and more intricate operations.

However, despite these limitations, the proposed IDMS represents a significant advancement over conventional logistic management methods. Traditional methods, often manual and dependent on human decision-making, lack the speed, precision, and adaptability that our system offers [35]. The implementation of the Swarm-Neural Network allows for better accuracy in forecasting, superior route optimization, and efficient inventory management, thereby enhancing overall operational efficiency.

Furthermore, the IDMS, being a machine learning-driven model, is capable of continuous learning and improvement, a distinct advantage over traditional systems. With time, as the system processes more data, it can refine its algorithms, enhance its predictive accuracy, and make more informed and effective decisions. Such an evolving capability of the IDMS presents a notable edge over static traditional systems.

Looking towards the future, there are several promising avenues to explore. One key perspective is to examine the integration of other AI techniques, such as reinforcement learning or deep learning, within the LBPSMS to complement the Swarm-Neural Network model. These techniques could further enhance the system's decision-making capabilities and adaptability. Further, future research could delve into the integration of the IDMS within a broader supply chain management framework, moving beyond logistics to explore applications in procurement, vendor management, or even customer relationship management.

In conclusion, our research demonstrated the significant potential of an IDMS in LBPSMS using ML models. The system's superior performance, adaptability, and decision-making capabilities highlight the transformative potential of integrating AI and ML in logistics management. While there are areas that require further research and challenges that need to be addressed, the advancements made in this study provide a solid foundation for future efforts in this direction. It is hoped that this research will catalyze further development in this field, contributing to the continuous evolution and improvement of logistics management systems. The next steps lie in expanding the scope of this research, delving into unexplored areas, and driving forward the frontier of ML-enabled logistics. The potential for an increasingly intelligent, efficient, and dynamic logistics sector is exciting, and we look forward to the continued progress in this field.

## VI. CONCLUSION

This research has comprehensively explored the development and implementation of an Intellectual Decision Making System (IDMS) within a Logistic Business Process Management System, harnessing the power of Machine Learning models. Our study has elucidated the transformative potential of integrating ML into the logistics sector, demonstrating substantial advancements in efficiency, decision-making, and operational optimization.

The proposed IDMS leverages swarm-neural network models, incorporating the strengths of swarm intelligence and neural networks to provide enhanced decision-making capabilities. Our system's exemplary performance in various areas of logistics, including demand forecasting, route optimization, and inventory management, reflects the efficacy and versatility of ML-enabled logistics systems. Particularly noteworthy was the system's adaptability under dynamic conditions, which underscores the value of ML models in navigating the complex, fluctuating landscapes of logistics operations.

However, our research is not without its limitations. While the IDMS demonstrated promising results, its application to other logistics areas and its scalability to larger operations remains untested. Additionally, the transition towards AI and ML in logistics poses challenges related to data privacy, computational resource demands, and the need for skilled personnel. These potential hurdles underline the need for strategic planning and comprehensive strategies in adopting AI-enabled logistics systems.

In closing, this research offers significant contributions to the field of ML-enabled logistics management, presenting a robust, adaptable, and efficient IDMS for LBPSMS. Our study not only confirms the advantages of ML in logistics but also paves the way for further exploration and development in this area. It is our hope that this research serves as a catalyst for future studies, fostering continuous evolution and innovation in logistics management. As the landscape of logistics continues to transform, we anticipate the continued growth and potential of AI and ML in redefining and enhancing the sector.

REFERENCES

[1] Zhou, L., Jiang, Z., Geng, N., Niu, Y., Cui, F., Liu, K., & Qi, N. (2022). Production and operations management for intelligent manufacturing: A systematic literature review. International Journal of Production Research, 60(2), 808-846.

[2] Yin, L., Zhong, R. R., & Wang, J. (2023). Ontology based package design in fresh E-Commerce logistics. Expert Systems with Applications, 212, 118783.

[3] Lei, N. (2022). Intelligent logistics scheduling model and algorithm based on Internet of Things technology. Alexandria Engineering Journal, 61(1), 893-903.

[4] Altayeva, A., Omarov, B., Suleimenov, Z., & Im Cho, Y. (2017, June). Application of multi-agent control systems in energy-efficient intelligent building. In 2017 Joint 17th World Congress of International Fuzzy Systems Association and 9th International Conference on Soft Computing and Intelligent Systems (IFSA-SCIS) (pp. 1-5). IEEE.

[5] UmaMaheswaran, S. K., Prasad, G., Omarov, B., Abdul-Zahra, D. S., Vashistha, P., Pant, B., & Kaliyaperumal, K. (2022). Major challenges and future approaches in the employment of blockchain and machine learning techniques in the health and medicine. Security and Communication Networks, 2022.

[6] Xu, X., & He, Y. (2022). Blockchain application in modern logistics information sharing: A review and case study analysis. Production Planning & Control, 1-15.

[7] Czvetkó, T., Kummer, A., Ruppert, T., & Abonyi, J. (2022). Data-driven business process management-based development of Industry 4.0 solutions. CIRP journal of manufacturing science and technology, 36, 117-132.

[8] Omarov, B., & Altayeva, A. (2018, January). Towards intelligent IoT smart city platform based on OneM2M guideline: smart grid case study. In 2018 IEEE International Conference on Big Data and Smart Computing (BigComp) (pp. 701-704). IEEE.

[9] Chen, Y., Li, R., & Song, T. (2023). Does TMT internationalization promote corporate digital transformation? A study based on the cognitive process mechanism. Business Process Management Journal, (ahead-of-print).

[10] Choi, T. M., Kumar, S., Yue, X., & Chan, H. L. (2022). Disruptive technologies and operations management in the Industry 4.0 era and beyond. Production and Operations Management, 31(1), 9-31.

[11] Malik, R., & Rybkowska, K. (2023). Business Processes Powered by Big Data: Current Issues and New Research Directions. In Big Data and Decision-Making: Applications and Uses in the Public and Private Sector (pp. 145-161). Emerald Publishing Limited.

[12] Liu, C., Feng, Y., Lin, D., Wu, L., & Guo, M. (2020). Iot based laundry services: an application of big data analytics, intelligent logistics management, and machine learning techniques. International Journal of Production Research, 58(17), 5113-5131.

[13] Shen, Z. M., & Sun, Y. (2023). Strengthening supply chain resilience during COVID - 19: A case study of JD. com. Journal of Operations Management, 69(3), 359-383.

[14] Helo, P., & Hao, Y. (2022). Artificial intelligence in operations management and supply chain management: An exploratory case study. Production Planning & Control, 33(16), 1573-1590.

[15] Sgarbossa, F., Grosse, E. H., Neumann, W. P., Battini, D., & Glock, C. H. (2020). Human factors in production and logistics systems of the future. Annual Reviews in Control, 49, 295-305.

[16] Dolgui, A., & Ivanov, D. (2022). 5G in digital supply chain and operations management: Fostering flexibility, end-to-end connectivity and real-time visibility through internet-of-everything. International Journal of Production Research, 60(2), 442-451.

[17] Gayialis, S. P., Kechagias, E. P., Konstantakopoulos, G. D., & Papadopoulos, G. A. (2022). A Predictive Maintenance System for Reverse Supply Chain Operations. Logistics, 6(1), 4.

[18] Tian, G., Lu, W., Zhang, X., Zhan, M., Dulebenets, M. A., Aleksandrov, A., ... & Ivanov, M. (2023). A survey of multi-criteria decision-making techniques for green logistics and low-carbon transportation systems. Environmental Science and Pollution Research, 1-23.

[19] Bai, L., Bai, J., & An, M. (2022). A methodology for strategy-oriented project portfolio selection taking dynamic synergy into considerations. Alexandria Engineering Journal, 61(8), 6357-6369.

[20] Haleem, A., Javaid, M., Singh, R. P., Suman, R., & Khan, S. (2023). Management 4.0: Concept, applications and advancements. Sustainable Operations and Computers, 4, 10-21.

[21] Altayeva, A. B., Omarov, B. S., Aitmagambetov, A. Z., Kendzhaeva, B. B., & Burkitbayeva, M. A. (2014). Modeling and exploring base station characteristics of LTE mobile networks. Life Science Journal, 11(6), 227-233.

[22] Ren, S. (2022). Optimization of enterprise financial management and decision-making systems based on big data. Journal of Mathematics, 2022, 1-11.

[23] Chen, X., Chen, R., & Yang, C. (2022). Research to key success factors of intelligent logistics based on IoT technology. The Journal of Supercomputing, 78(3), 3905-3939.

[24] A. Altayeva, B. Omarov, H.C. Jeong, Y.I. Cho. Multi-step face recognition for improving face detection and recognition rate. Far East Journal of Electronics and Communications 16(3), pp. 471-491.

[25] Lu, M., & Wudhikarn, R. (2022, January). Using the best-worst method to develop intellectual capital indicators in financial service company. In 2022 Joint International Conference on Digital Arts, Media and Technology with ECTI Northern Section Conference on Electrical, Electronics, Computer and Telecommunications Engineering (ECTI DAMT & NCON) (pp. 81-86). IEEE.

[26] Lăzăroiu, G., Andronie, M., Iatagan, M., Geamănu, M., Ştefănescu, R., & Dijmărescu, I. (2022). Deep Learning-Assisted Smart Process Planning, Robotic Wireless Sensor Networks, and Geospatial Big Data Management Algorithms in the Internet of Manufacturing Things. ISPRS International Journal of Geo-Information, 11(5), 277.

[27] Alsudani, M. Q., Jaber, M. M., Ali, M. H., Abd, S. K., Alkhayyat, A., Kareem, Z. H., & Mohhan, A. R. (2023). Smart logistics with IoT-based enterprise management system using global manufacturing. Journal of Combinatorial Optimization, 45(2), 57.

[28] Sun, X., Yu, H., & Solvang, W. D. (2022). Towards the smart and sustainable transformation of Reverse Logistics 4.0: a conceptualization and research agenda. Environmental Science and Pollution Research, 29(46), 69275-69293.

[29] Li, C., Chen, Y., & Shang, Y. (2022). A review of industrial big data for decision making in intelligent manufacturing. Engineering Science and Technology, an International Journal, 29, 101021.

[30] Pereira, A. M., Moura, J. A. B., Costa, E. D. B., Vieira, T., Landim, A. R., Bazaki, E., & Wanick, V. (2022). Customer models for artificial intelligence-based decision support in fashion online retail supply chains. Decision Support Systems, 158, 113795.

[31] Aboamer, M. A., Sikkandar, M. Y., Gupta, S., Vives, L., Joshi, K., Omarov, B., & Singh, S. K. (2022). An investigation in analyzing the food quality well-being for lung cancer using blockchain through cnn. Journal of Food Quality, 2022.

[32] Cherchata, A., Popovychenko, I., Andrusiv, U., Gryn, V., Shevchenko, N., & Shkuropatskyi, O. (2022). Innovations in logistics management as a direction for improving the logistics activities of enterprises. Management Systems in Production Engineering, 30(1), 9-17.

[33] Choi, T. M., & Siqin, T. (2022). Blockchain in logistics and production from Blockchain 1.0 to Blockchain 5.0: An intra-inter-organizational framework. Transportation Research Part E: Logistics and Transportation Review, 160, 102653.

[34] Ren, S., Choi, T. M., Lee, K. M., & Lin, L. (2020). Intelligent service capacity allocation for cross-border-E-commerce related third-party-forwarding logistics operations: A deep learning approach. Transportation Research Part E: Logistics and Transportation Review, 134, 101834.

[35] Rejeb, A., Simske, S., Rejeb, K., Treiblmaier, H., & Zailani, S. (2020). Internet of Things research in supply chain management and logistics: A bibliometric analysis. Internet of Things, 12, 10031.

# Construction of Short-Term Traffic Flow Prediction Model Based on IoT and Deep Learning Algorithms

Xiaowei Sun[1], Huili Dou[2]*

Zhejiang Institute of Communications, Hangzhou 310012, China[1]

Institute of Rail Transit, Zhejiang Institute of Communications, Hangzhou 310012, China[2]

*Abstract*—On a global scale, traffic problems are an essential factor affecting urban operations, particularly challenging the frequent occurrence of traffic congestion and accidents. The solution to the problem requires real-time and accurate prediction of traffic flow. This article mainly explores the application of the Internet of Things and deep learning in traffic flow prediction, aiming to solve the problem where existing methods cannot meet the requirements of real-time and accuracy. IoT devices, such as road sensors and in-vehicle GPS devices, which provides rich information for traffic flow prediction. With the ability of deep learning, it can not only learn and abstract a large amount of complex traffic data but also handle traffic flow prediction tasks in various complex situations. During the model construction process, the complexity of the road network was fully considered, practical algorithms were designed to fuse multi-source data, and the structure of the model was optimized to meet the needs of real-time prediction. The experimental results show that the absolute error of the test results is generally less than 6km/h, which can better reflect the traffic speed of the road section in the future.

*Keywords—Internet of things; deep learning algorithm; short term traffic flow; prediction model*

## I. INTRODUCTION

With increase in serious urban transportation problems, seriously affecting the functional operation of cities and the quality of life of citizens. Among them, the frequent occurrence of traffic congestion and accidents has become a common problem worldwide [1, 2]. Traffic flow prediction is crucial for urban traffic management [3, 4], as it can accurately and effectively predict traffic flow. Traffic management departments can arrange traffic police forces in advance, dispatch traffic lights reasonably, and effectively guide vehicles based on the prediction results, thereby alleviating traffic congestion and improving urban road capacity [5]. For example, when it is predicted that the traffic flow in a specific area will significantly increase, traffic control or evacuation work can be carried out in advance to avoid traffic congestion. For drivers and passengers, knowing the traffic flow situation in the future in advance can provide important references for their travel decisions, reduce waiting time, and improve travel efficiency [6]. Therefore, traffic flow prediction is also significant for the research and development. In areas such as autonomous driving, traffic signal optimization, and travel recommendation, accurate traffic flow prediction results are all needed [7].

In the information age of the 21st century, deep learning has seen significant growth in many fields. Traffic flow

prediction, as a critical challenge, is gradually benefiting from these two technologies. The Internet of Things, also known as the extension of the Internet, connects various objects in the physical world, enabling them to collect and exchange data. The Internet of Things technology has played a considerable role in traffic flow prediction [8]. The use of IOT makes it possible to obtain large amounts of data, which greatly improves the accuracy and real-time nature of data sources. Deep learning, as an artificial intelligence algorithm, has also played an enormous role in traffic flow prediction. Deep learning can learn and understand a large amount of complex data and abstract valuable features [9, 10]. It has been successfully applied to traffic flow prediction. The application of deep learning enables traffic flow prediction models to understand and process complex traffic data, thereby improving the accuracy of predictions. Although the Internet of Things and deep learning have been successfully applied in traffic flow prediction, their potential still needs to be fully explored [11,12]. In terms of deep learning, how to design more effective models to handle more complex situations (such as traffic congestion, accidents, etc.) is also an important research direction. Therefore, constructing traffic flow prediction models based on the Internet of Things and deep learning is an essential direction in current traffic research.

However, although the Internet of Things and deep learning have been applied in traffic flow prediction, they still need to overcome many challenges. Data quality issues, such as sensor failures, network transmission issues, incomplete data, etc., can all affect the availability and accuracy of data [13, 14]. In terms of model design and optimization, how to design and optimize the model based on specific traffic flow prediction tasks and how to improve the interpretability of the model are all issues that need to be addressed. Real-time prediction problems require the model to have efficient data processing and computational capabilities to meet the needs of real-time prediction. The complexity of road networks and how to effectively integrate various factors into the model is a challenging issue. For the fusion of multi-source data, in addition to traffic flow data, weather data, social event data, social media data, etc., can also be utilized. How to effectively integrate these diverse data and improve prediction accuracy is a new challenge. Therefore, the significance of studying the short-term traffic flow prediction model based on the Internet of Things and deep learning algorithm is that it can significantly improve the efficiency of urban traffic management, reduce traffic congestion, and reduce the incidence of traffic accidents. Through real-time and accurate traffic flow prediction, traffic management departments can

allocate resources reasonably, optimize traffic signal scheduling, and guide vehicles to drive effectively. At the same time, this research has a significant impact on the development of intelligent transportation systems, especially in the fields of autonomous driving and traffic signal optimization.

## II. APPLICATION OF IoT AND DEEP LEARNING

### A. *Deep Learning Algorithm*

The architecture of the Temporal Convolutional Network (TCN) is designed based on the characteristics of the latest convolutional system used for sequential data. It takes a time series as input and models the temporal correlation in each temporal data. Unlike the traditional Recurrent Neural Network (RNN) which recourses along the time axis of a sequence and introduces a large number of learning parameters, making the model difficult to optimize. TCN combines simplicity, autoregressive prediction, and very long memory without the recursive mechanism of RNN, which facilitates parallelization, as shown in Fig. 1. Therefore, TCN can effectively combine computational advantages with representation capabilities to achieve efficient and good predictive performance [15]. Due to the above advantages, TCN can be well applied to analyse data with strict order and is widely used for predicting various scenarios. Therefore, TCN is suitable for modelling and analysing time-series data sensors monitor. By capturing simple patterns in sensor time series and generating more complex patterns in higher-level layers, TCN can better extract temporal features [16].

Graph Convolutional Neural Networks (GCN) are feature extractors designed based on graph data [17]. The essence of GCN is to apply convolution to graph neural networks, which can flexibly extract structural information of graph data and reduce computational complexity. Due to the good complementary relationship between node attribute information and structural information in graph data, GCN can use the network layer to simultaneously learn the data structure and attribute information in the graph and use the two to represent the relationship between nodes [18, 19].

### B. *Application of Deep Learning in Traffic Flow Prediction*

*1) Reactive control of short-term traffic flow:* The timing control method calculates the timing scheme based on historical traffic flow data and predetermined optimization objectives. The biggest drawback of this scheme is that it cannot adapt to the dynamic changes in traffic flow, resulting in limited control effectiveness. The reactive traffic signal control method adjusts the signal timing strategy based on existing traffic flow characteristics to improve control effectiveness without considering the impact of traffic flow prediction on control effectiveness [20, 21]. In recent years, green ratio optimization, phase difference optimization, mathematical programming, multi-objective optimization, dynamic programming and other methods have emerged in the field of reactive control [22]. Based on the analysis and research of the delay law of the vehicles at the intersection, an optimization model of phase difference adjustment of the wire control system is established. On the one hand, due to the lack of traffic flow forecasting mechanism, the control method does not consider the potential impact of the current control scheme on the future traffic conditions, so the control effect is limited. On the one hand, the control method lacks a traffic flow prediction mechanism. It needs to consider the potential impact of the current control plan on future traffic conditions, resulting in limited control effectiveness. On the other hand, most of these existing control methods adopt a single-machine computing environment, which cannot meet the real-time requirements of traffic optimization and control in the context of big data [23, 24].



Fig. 1. Deep learning applied to short-term traffic models.

Fig. 2. Construction process of short-term traffic network image samples.

Similar to facial muscle movements forming different expressions, the spatiotemporal evolution of short-term traffic flow in the road network constitutes different forms of traffic. From a spatial perspective, short-term traffic network flows are interdependent and interrelated. Congestion on a road section may affect nearby or even further road sections. At the same time, in terms of time dimension, some similar but randomly fluctuating traffic flow characteristics will repeatedly appear on a road segment. Therefore, short-term traffic network flow feature learning needs to comprehensively consider the temporal and spatial characteristics, as well as periodic repeatability characteristics. As shown in Fig. 2, the construction process of short-term traffic network image samples is presented.

*2) Short-term traffic flow model predictive control:* Based on the model predictive control (MPC) framework, the predictive model predicts the future traffic dynamics, and the potential control performance of the candidate scheme is calculated [25, 26]. In the current control cycle, the first element of the optimization sequence is applied to the traffic system model and restarts the next round of the rolling optimization process based on the feedback traffic status and prediction model. By comprehensively considering the

cumulative impact trends, visionary control decisions are generated. Therefore, in urban expressway traffic control, MPC can synergistically adjust the traffic flow at the ramp entrance and exit. In highway traffic control, MPC can coordinate and solve problems such as speed restrictions, lane allocation, and release time of vehicle queues on on-ramps. The rolling time domain method of MPC can further plan the path selection problem of multiple travellers. Abstract boundary control and path guidance as economic MPC problems is to improve the mobility of urban networks. Furthermore, in the MPC framework, macro traffic flow and exhaust emission models are introduced to reduce the probability of traffic congestion and reduce pollution emissions. However, to apply the MPC control strategy to more complex traffic network systems, the contradiction between the time required optimizing the objective function in the prediction time domain and the real-time performance of online control still needs to be solved urgently [27].

In order to reduce the solving time of the MPC objective function, the entire road network is decomposed into several regional subnets to accelerate the calculation process. In addition, by parameterizing the macro traffic prediction model,

it reduces the time of online computation in the rolling time domain. Based on the improved macro traffic flow model, and quadratic programming provide a solution for improving the real-time performance of MPC in traffic flow control. However, they reduce the calculation time for solving the objective function in the prediction time domain. They can only roughly describe macroscopic traffic flow phenomena such as traffic density, traffic flow, occupancy rate, etc., which, to some extent, reduces the control effect.

*3) Deep reinforcement learning control for short-term traffic flow:* Deep reinforcement learning is a feedback-based iterative learning method based on a deep learning evaluation mechanism. Deep learning involves constructing a hierarchical neural network that simulates human brain thinking. The development of deep learning to this day mainly includes CNN, deep belief networks DBN, DSAEs, and LSTM, each with its advantages and applicability. However, these deep learning methods focus on the learning of traffic flow characteristics at the segment level. When traffic control rises to the regional road network, the best control timing may be missed due to the inability to obtain the interconnectivity between segment traffic flows [28, 29].

The training of deep learning in the context of big data takes several days or even weeks, so reducing the time cost of model training while ensuring training accuracy has become a hot topic in academic and industrial research [30]. Distributed deep learning is a powerful tool to accelerate deep neural network training. By introducing predictive hierarchical caching strategy in distributed training, it can improve the cache hit rate of data, shorten synchronization time and network blocking times. Secondly, through the sparse gradient compression mechanism of entropy, the propagation gradient threshold can be determined dynamically and the data volume of the propagation gradient can be compressed to reduce the communication load. By quantifying the performance differences of each node and dynamically allocating the training batches of each node, the time of each iteration between nodes is approximately consistent, thereby improving the impact of gradient obsolescence on convergence in asynchronous parallel optimization [31].

*4) Distributed parallel processing of traffic flow big data:* The rapid development of the Internet of Things and artificial intelligence has provided strong support for the interconnection of vehicles, pedestrians, traffic lights, roadside equipment, and traffic management centers. It is necessary to establish new theories and methods for traffic network flow adaptive control and achieve the next generation of data-driven intelligent transportation systems (ITS). In the context of the Internet of Vehicles (IoV), the traffic flow data collected by multi-source heterogeneous sensors is rapidly increasing, and the era of big data in transportation has arrived. Cloud computing uses a universal computing model to deploy computing tasks to a computing resource pool, allowing users to transparently access computing resources, storage space, and information services according to their needs. It is one of the most effective methods for processing big data. It has the self-maintenance and management function of virtual computing resources and can dynamically acquire or release computing resources to adapt to dynamic application workloads.

In a traffic control system, if all raw data is sent to a remote traffic control center for processing and analysis using cloud computing, it requires extremely high network bandwidth. In addition, when optimization decisions are returned from cloud computing centers to traffic signal controllers, local traffic dynamics may have undergone significant changes. This poses hidden dangers to the safety and real-time performance of traffic control. Edge computing is an expansion of cloud computing architecture, pushing some computing intelligence, data processing, storage, and services from the cloud to the network's edge. It enables analysis and processing to occur on the side of the data source, avoiding response delays or data security risks caused by long-distance, high-capacity data communication as much as possible.

## III. CONSTRUCTION OF SHORT-TERM TRAFFIC FLOW PREDICTION MODEL BASED ON IOT AND DEEP LEARNING ALGORITHMS

### A. Overall Architecture

As shown in Fig. 3, the model predictive control architecture is deployed on a cloud computing platform to collaboratively control the signal timing strategies of various intersections from a global perspective in order to improve the traffic capacity of vehicles in the road network and alleviate traffic congestion. Establish information channels between transportation networks and cloud computing through communication technologies like the Internet and 5G. The location, speed, and intersection status of vehicles in the transportation network are collected through multi-source sensors and then uploaded to the cloud control center. The nonanalytical prediction model of the cloud control center predicts the trend of traffic flow changes in the future based on the traffic status collected at the current time, pre-set control requirements, and control the sequences generated by optimization algorithms and provides an evaluation of the cumulative control performance of the control sequence in the future. Using the distributed computing of cloud computing, multiple computing nodes participate in the calculation to accelerate the optimal control sequence solution. In the current control cycle, the first strategy in the optimal control sequence is selected and applied to the traffic network flow system. The rolling time domain method is used to continuously implement this process and effectively control the traffic flow of the road network.

Fig. 3. Overall architecture diagram.

## B. Short Term Traffic Flow Simulation Modelling

In urban road networks, vehicles face complex road conditions on their driving routes, as shown in Fig. 4.

*1) Maximum speed limit:* The maximum allowable driving speed depends on the road infrastructure and equipment. When the road conditions and driving equipment performance are good, the maximum driving speed can increase correspondingly. Otherwise, it is necessary to reduce the maximum driving speed to meet actual needs.

*2) District-specific speed limits:* The urban road network has many special areas, such as hospitals, schools, military administration areas, and signalized intersections, where vehicles need to slow down appropriately.

*3) Temporary speed limit:* When encountering sudden situations such as roadbed maintenance, abnormal weather, and traffic accidents, the relevant traffic management department will issue a temporary speed limit notice, and vehicles should slow down in advance and pass slowly when driving to the section. $V_{lim}(x)$ dynamically divides the road into a series of sections with different speed restrictions.



Fig. 4. Spatiotemporal constraints on vehicle motion.

## C. Establishment of Short-term Traffic Flow Model Prediction Model

Decompose the complete dataset D into R parts, with each part represented by Dr. (r=1, 2,.. R). Therefore, dataset D can be represented as a set of subsets of data, Dr., as shown in Formula (1):

$$D = U_{r=1}^{R} D^{r} = U_{r=1}^{R} U_{n=1}^{N} D_{n}^{r} \qquad (1)$$

Among them, Nr represents the size of the r-th data subset, and Dr n represents the nth data on the r-th data subset.

The objective function of parallel training of CNN-LSTM model is shown in Formula (2):

$$J = \min \frac{1}{N} \sum_{r=1}^{R} j^{r} = \min \frac{1}{N} \sum_{r=1}^{R} \sum_{n=1}^{N} j_{n}^{r} \qquad (2)$$

order $j_n^r = J_n^r = \|y_n^r - \hat{y}_n^r\|^2 / 2M$, Here $y_n^r$ and $\hat{y}_n^r$ are the observation and prediction vectors for the n-th sample

in the r-th dataset, respectively, The training of the CNN-LSTM model based on the complete dataset D is to minimize the objective function described by the formula, thereby obtaining ideal weights and biases, known as global learning parameters. Furthermore, The weights and biases trained by minimizing local objective function Jr on data subset Dr are called local learning parameters. For parallel feature forward learning processes, the output values of different types of network layers are synchronously calculated in parallel based on corresponding data subsets. At time t, denoted by $a_{n,j,c}^{r,l}(t)$, the CNN layer extracts local feature values based on the data subset Dr. The calculation Formula (3) is as follows:

$$a_{n,j,c}^{r,l}(t) = \sigma\left(\sum_{i=1}^{N_c^{l-1}} a_{n,i,c}^{r,l-1}(t) * \omega_{j,i,c}^{r,l} + b_{j,c}^{r,l}\right) \qquad (3)$$

Among them, c represents the convolutional layer of the CNN-LSTM model:



Fig. 5. Structural diagram of LSTM.

As shown in Fig. 5, in the LSTM module, the output values of forgetting gates (such as formulas), input gates (such as formulas), cellular states (such as formulas), output gates (such as formulas), and implicit states (such as formulas) at time t are represented by $\mathrm{a}^{r,l}_{n,j,lm,fg}(t)$ , $a^{r,l}_{n,j,lm,ig}(t)$ , $a^{r,l}_{n,j,lm,cg}(t)$ , $a^{r,l}_{n,j,lm,og}(t)$ , and, respectively. The calculation formulas for these locally activated feature values are shown in (4) - (9):

$$a^{r,l}_{n,j,lm,fg}(t) = \sigma(\sum_{i=1}^{N^l_{hg}} w^{r,l}_{j,i,lm,fg}(t-1) + \sum_{i=N^l_{hg}+1}^{N^l_{hg}+N^{l-1}_c} w^{r,l}_{j,i,lm,fg} a^{r,l-1}_{n,i,c}(t) + b^{r,l}_{j,lm,fg})$$
(4)

$$a^{r,l}_{n,j,lm,ig}(t) = \sigma(\sum_{i=1}^{N^l_{hg}} w^{r,l}_{j,i,lm,hg} a^{r,l}_{n,i,lm,hg}(t-1) + \sum_{i=N^l_{hg}+1}^{N^l_{hg}+N^{l-1}_c} w^{r,l}_{j,i,lm,ig} a^{r,l-1}_{n,i,c}(t) + b^{r,l}_{j,lm,ig})$$
(5)

$$a^{r,l}_{n,j,lm,cg}(t) = a^{r,l}_{j,i,lm,fg}(t) a^{r,l}_{n,i,lm,cg}(t-1) + a^{r,l}_{n,i,lm,ig}(t) \tilde{a}^{r,l}_{n,i,lm,cg}(t)$$
(6)

$$\tilde{a}^{r,l}_{n,j,lm,cg}(t) = \tanh(\sum_{i=1}^{N^l_{hg}} w^{r,l}_{j,i,lm,cg} a^{r,l}_{n,i,lm,hg}(t-1) + \sum_{i=N^l_{hg}+1}^{N^l_{hg}+N^{l-1}_c} w^{r,l}_{j,i,lm,cg} a^{r,l-1}_{n,i,c}(t) + b^{r,l}_{j,lm,cg})$$
(7)

$$a^{r,l}_{n,j,lm,og}(t) = \sigma(\sum_{i=1}^{N^l_{hg}} w^{r,l}_{j,i,lm,og} a^{r,l}_{n,i,lm,hg}(t-1) + \sum_{i=N^l_{hg}+1}^{N^l_{hg}+N^{l-1}_c} w^{r,l}_{j,i,lm,og} a^{r,l-1}_{n,i,c}(t) + b^{r,l}_{j,lm,og})$$
(8)

$$a^{r,l}_{n,j,lm,hg}(t) = a^{r,l}_{n,i,lm,og}(t) \tanh(a^{r,l}_{n,i,lm,cg}(t))$$
(9)

Let $a^{r,l}_{n,j,f}(t)$ represent the network output of the fully connected layer at time t, and the calculation Formula is shown in (10):

$$a^{r,l}_{n,j,f}(t) = \sigma(\sum_{i=1}^{N^{l-1}_f} w^{r,l-1}_{j,i,f}(t) + b^{r,l}_{j,f})$$
(10)

Among them, f represents the fully connected layer. When layer l is a fully connected layer, i represents the i-th input neuron, j represents the j-th output neuron. $w^{r,l-1}_{j,i,f}(t)$ and $b^{r,l}_{j,f}$ represent the weight and bias of layer l (which is a fully connected layer) on the r-th data subset, respectively, and $N^{l-1}_f$ represents the number of neurons in the previous layer; σ(·) represents the RELU activation function.

Based on the classical gradient descent criterion, the relationship between global learning parameters and local learning parameters in the parallel error backpropagation process is derived layer by layer. The calculation formulas for updating the global learning parameters $w^l_{j,i,f}(t)$ and $b^l_{j,f}(t)$ in the fully connected layer at time step t are shown in Formulas (11) - (12):

$$w^l_{j,i,f}(t) = w^l_{j,i,f}(t-1) - \eta \frac{\partial J}{\partial w^l_{j,i,f}} = \frac{1}{R}\sum_{r=1}^{R} w^l_{j,i,f}(t)$$
(11)

$$b^l_{j,f}(t) = b^l_{j,f}(t-1) - \eta \frac{\partial J}{\partial b^l_{j,f}} = \frac{1}{R}\sum_{r=1}^{R} b^l_{j,f}(t)$$
(12)

Among them, the local weights and biases in the fully connected layer are calculated as shown in Formulas (13) - (14):

$$w^{r,l}_{j,i,f}(t) = w^l_{j,i,f}(t-1) - \frac{\eta}{N}\sum_{n=1}^{N^r} R\delta^{r,l}_{n,j,f} a^{r,l-1}_{n,i,f}$$
(13)

$$b^{r,l}_{j,f}(t) = b^l_{j,f}(t-1) - \frac{\eta}{N}\sum_{n=1}^{N^r} R\delta^{r,l}_{n,j,f}$$
(14)

The global adaptive learning rate of parallel training of the CNN-LSTM model can be obtained by calculating the local gradient sum, as shown in Formulas (15) - (16):

$$\eta(t) = \frac{l_r}{\mu + \sqrt{G(t)}}$$
(15)

$$G(t) = \rho G(t-1) + (1-\rho)g(t)$$
(16)

Among them, G(t) represents the sum of squares of gradients with attenuation factors; $\rho$ Similar to the attenuation factor in the momentum gradient descent method, it represents the impact of past gradients on current parameter updates, typically taking a value of 0.9; Lr represents the basic learning rate; μ It is a minimal constant that prevents the denominator from being zero. The adaptive learning rate has advantages over the traditional fixed learning rate, because it can adjust the learning rate according to the gradient of the parameter itself, so as to achieve better convergence effect. Regardless of whether the training data set D is evenly divided or inevenly divided, the adaptive learning rate ensures that the convergence results are almost identical to the serial training method. Therefore, the parallel training theory of the CNN-LSTM model ensures that the global learning features of the large dataset can be obtained from the parallel learning of each decomposed data subset.

## IV. MODEL EXPERIMENT AND RESULT ANALYSIS

Fig. 6 shows the trend of CP curves for MAE indicators using different prediction methods. For the prediction tasks of traffic network flow in 5min, 15min, 30min, and 60min, the CNN-LSTM prediction method has 100%, 85.71%, 85.71%, and 71.4% of MAE errors controlled within 20 for expressways, respectively. Compared to other prediction methods, the CP curve of the CNN-LSTM prediction method is always located at the top left of the graph in different traffic network flow prediction tasks, indicating that the CNN-LSTM method has more advantages in improving the accuracy of traffic network flow prediction. The universal ability measures the adaptability of the CNN LSTM model to prediction tasks in

different traffic scenarios. As can be seen from the figure, MAE prediction error of CNN-LSTM method in different traffic network flow prediction tasks is mainly kept between 12.31 and 19.05. This shows that CNN-LSTM model has competitive adaptability in various prediction scenarios under the same prediction time domain. However, the MAE indices of DTR and SVR methods fluctuate greatly under different prediction tasks. These two prediction models are susceptible to differences in the prediction time domain, so their predictive universal ability could be better in different traffic scenarios. CNN-LSTM is more stable than other methods in terms of universal prediction ability. Through comparison, it was further found that the prediction accuracy of the CNN-LSTM method in larger prediction time domain tasks is lower than that in smaller prediction time domain tasks. This is because under the same training cycle, as the prediction time domain increases, the difficulty of predicting the future multi-step traffic dynamic evolution trend increases. In the future, the deep learning model structure will be improved to enhance the feature

extraction ability of traffic network flow in larger prediction time domains.

MPC online optimization is carried out rolling, and the end of the current control cycle starts the subsequent predictive time domain optimization. Therefore, short-term traffic network flow predictive control based on the rolling time domain includes multiple predictive time domain optimization processes. Fig. 7 shows the computational efficiency of predicting time domain optimization for each control cycle. The MPC control scheme based on Spark cloud parallel optimization takes much less time at each control time step than the single machine serial optimization MPC control scheme. Especially for single-machine computing environments, all chromosomes sequentially call the traffic network flow prediction model circularly to obtain evaluation values for control effectiveness. This calculation method requires a high computational time cost for non-analytical micro prediction models.



Fig. 6. Cumulative distribution traffic flow prediction.

Fig. 7. Computational efficiency of predicting time-domain optimization for each control cycle.

On the contrary, on the Spark cloud, multiple chromosomes are divided into several subpopulations and distributed to various worker nodes. The more nodes there are, the fewer chromosomes are assigned to each worker node, resulting in a more minor computational task. The chromosomes on all worker nodes call the traffic network flow prediction model to obtain the control effect evaluation values of chromosomes in a parallel manner, thereby reducing optimization time. As shown in Fig. 7, at the beginning of the simulation, the acceleration ratio of parallel optimization based on Spark cloud is lower than that of single-machine serial optimization. With the deepening of the simulation, the efficiency of the late parallel computation is improved significantly and remains stable. This is because Spark initially takes some time to load the data. However, after data is cached to the memory, Spark can directly obtain data from the memory when it is invoked again, which improves computing efficiency.

From Fig. 8, it can be observed that the prediction results of the nonparametric regression algorithm are better than those of the BP neural network method. However, the improved algorithm in this paper has a particular improvement in prediction accuracy compared to the basic nonparametric regression algorithm, and the absolute error of the prediction results is generally below 6km/h, which can better reflect the future traffic speed situation of the road segment. In summary, it can be seen that the improved prediction method in this paper shows good prediction ability in the overall traffic speed fitting and specific prediction results. Compared with the general algorithm, this method has obvious improvement in prediction accuracy.



Fig. 8. Error distribution over time series.

## V. CONCLUSIONS

This article conducts in-depth research on short-term traffic flow prediction and constructs a prediction model based on the Internet of Things and deep learning algorithms. Through analysis of actual traffic data and model validation, we have drawn the following conclusions:

Many real-time traffic data, including vehicle sensors, traffic cameras, and traffic lights, has been obtained using

Internet of Things technology. These data provide comprehensive and accurate traffic status information, providing an essential foundation for short-term traffic flow prediction. By comparing the experimental results, it was found that deep learning algorithms have better performance and accuracy in traffic flow prediction tasks compared to traditional machine learning algorithms.

The prediction model proposed by this research institute based on the Internet of Things and deep learning algorithms provides strong support for actual traffic management and decision-making. This model can help traffic management departments better plan road resources, optimize traffic signal timing, and provide real-time traffic congestion information to drivers and traffic participants to improve traffic efficiency and reduce traffic congestion.

## ACKNOWLEDGMENT

## FUNDING

## REFERENCES

[1] Xia M. Traffic congestion index calculation based on BP neural network[J]. Advances in Computer, Signals and Systems, 2021, 5(1).

[2] Hassija V, Gupta V ,Garg S , et al. Traffic Jam Probability Estimation Based on Blockchain and Deep Neural Networks. IEEE Transactions on Intelligent Transportation Systems, 2020, PP(99).

[3] Yaqin Y, Yue X, Yuxuan Z, et al. Dynamic multi-graph neural network for traffic flow prediction incorporating traffic accidents. Expert Systems With Applications, 2023, 234.

[4] Xijun Z ,Jiwen L . Traffic flow prediction based on GRU-BP combined neural network. Journal of Physics: Conference Series, 2021, 1873(1).

[5] Guo Y, Lu L. Application of a Traffic Flow Prediction Model Based on Neural Network in Intelligent Vehicle Management. International Journal of Pattern Recognition and Artificial Intelligence, 2019, 33(3).

[6] Yi L, Mingsheng L,Yunchi X, et al. Traffic flow prediction model based on gated graph convolution with attention. Journal of Physics: Conference Series, 2023, 2493(1).

[7] Ameya A K, Shravan R, Ananya D , et al. Traffic flow prediction models – A review of deep learning techniques. Cogent Engineering, 2022, 9(1).

[8] Oreja M J, Gozalvez J. A Comprehensive Evaluation of Deep Learning-Based Techniques for Traffic Prediction. IEEE Access, 2020, 8.

[9] Yang L, Yaolun S,Yan Z, et al. WT-2DCNN: A convolutional neural network traffic flow prediction model based on wavelet reconstruction. Physica A: Statistical Mechanics and its Applications, 2022, 603.

[10] Ismaeel G A, Janardhanan K ,Sankar M , et al. Traffic Pattern Classification in Smart Cities Using Deep Recurrent Neural Network. Sustainability, 2023, 15(19).

[11] Qianqian Z, Nan C, Siwei L. FASTNN: A Deep Learning Approach for Traffic Flow Prediction Considering Spatiotemporal Features. Sensors, 2022, 22(18).

[12] Yuanmeng Z, Jie C, Hong Z, et al. A deep learning traffic flow prediction framework based on multi-channel graph convolution. Transportation Planning and Technology, 2021, 44(8).

[13] Zhao Z, Hao Y, Xianfeng Y. A Transfer Learning–Based LSTM for Traffic Flow Prediction with Missing Data. Journal of Transportation Engineering, Part A: Systems, 2023, 149(10).

[14] Bernardo G, José C, Helena A. A survey on traffic flow prediction and classification. Intelligent Systems with Applications, 2023, 20.

[15] Chen C ,Ziye L ,Shaohua W , et al. Traffic Flow Prediction Based on Deep Learning in Internet of Vehicles. IEEE Transactions On Intelligent Transportation Systems, 2021, 22(6).

[16] Jiang L, Luofeng J. Traffic Flow Prediction Method Based on Deep Learning. Journal of Physics: Conference Series, 2020, 1646 (1).

[17] Hong Z, Sunan K ,XiJun Z , et al. Dynamic Spatial–Temporal Convolutional Networks for Traffic Flow Forecasting. Transportation Research Record, 2023, 2677(9).

[18] Emerging Technologies; Research Conducted at Beijing Institute of Technology Has Updated Our Knowledge about Emerging Technologies (A hybrid deep learning based traffic flow prediction method and its understanding). Computers, Networks & Communications, 2018.

[19] Wu Y, Tan H, Qin L, et al. A hybrid deep learning based traffic flow prediction method and its understanding. Transportation Research Part C, 2018, 90.

[20] Yingya G, Yufei P, Run H, et al. Capturing spatial–temporal correlations with Attention based Graph Convolutional Network for network traffic prediction. Journal of Network and Computer Applications, 2023, 220.

[21] Siyuan F, Shuqing W, Junbo Z, et al. A macro–micro spatio-temporal neural network for traffic prediction. Transportation Research Part C, 2023, 156.

[22] Bo W, L. H V, Inhi K, et al. Distributional prediction of short-term traffic using neural networks. Engineering Applications of Artificial Intelligence, 2023, 126(PC).

[23] Rui H, Cuijuan Z, Yunpeng X, et al. Deep spatio-temporal 3D dilated dense neural network for traffic flow prediction. Expert Systems with Applications, 2024, 237(PA).

[24] Xian Y, Yin-Xin B, Quan S. STHSGCN: Spatial-temporal heterogeneous and synchronous graph convolution network for traffic flow prediction. Heliyon, 2023, 9(9).

[25] Robert J ,Young T K ,Emily G , et al. Tailoring Mission Effectiveness and Efficiency of a Ground Vehicle Using Exergy-Based Model Predictive Control (MPC). Energies, 2021, 14(19).

[26] Artificial Neural Network; Findings on Artificial Neural Network Discussed by Investigators at Ryerson University. Internet Networks & Communications, 2017.

[27] Sunday S O. The Application of Model Predictive Control (MPC) to Fast Systems such as Autonomous Ground Vehicles (AGV). IOSR Journal of Computer Engineering, 2014, 16(3).

[28] Yaqin Y, Yue X, Yuxuan Z, et al. Dynamic multi-graph neural network for traffic flow prediction incorporating traffic accidents. Expert Systems with Applications,2023,234.

[29] Yi X, Liangzhe H, Tongyu Z, et al. Generic Dynamic Graph Convolutional Network for traffic flow forecasting. Information Fusion, 2023, 100.

[30] Dongran Z, Jun L. Multi-view fusion neural network for traffic demand prediction. Information Sciences,2023,646.

[31] Xiaoxiao S, Xinfeng W, Boyi H, et al. Multidirectional short-term traffic volume prediction based on spatiotemporal networks. Applied Intelligence, 2023, 53(20).

# Deep Learning for Early Detection of Tomato Leaf Diseases: A ResNet-18 Approach for Sustainable Agriculture

Asha M S[1], Yogish H K[2]

Department of Computer Science and Engineering, Christ University (Deemed to be University), Karnataka, India[1]
Department of Information Science and Engineering, M S Ramaiah Institute of Technology, Karnataka, India[2]
Visvesvaraya Technological University, Belagavi, Karnataka, India-590018[1, 2]

*Abstract*—The paper explores the application of Convolutional Neural Networks (CNNs), specifically ResNet-18, in revolutionizing the identification of diseases in tomato crops. Facing threats from pathogens like Phytophthora infestans, timely disease detection is crucial for mitigating economic losses and ensuring food security. Traditionally, manual inspection and labour-intensive tests posed limitations, prompting a shift to CNNs for more efficient solutions. The study uses a well-organized dataset, employing data preprocessing techniques and ResNet-18 architecture. The model achieves remarkable results, with a 91% F1 score, indicating its proficiency in distinguishing healthy and unhealthy tomato leaves. Metrics such as accuracy, sensitivity, specificity, and a high AUC score on the ROC curve underscore the model's exceptional performance. The significance of this work lies in its practical applications for early disease detection in agriculture. The ResNet-18 model, with its high precision and specificity, presents a powerful tool for crop management, contributing to sustainable agriculture and global food security.

*Keywords—Convolution neural networks; tomato crop health; deep learning; binary classification; disease detection*

## I. INTRODUCTION

Tomato (Solanum Lycopersicon) holds an important place in agriculture, food and cooking worldwide. Known for its bright red color and its many uses in dishes such as Mediterranean pastas, Asian curries, and American tomato sauces, the tomato has become an important part of world cuisines [1]. In addition to its gastronomic importance, tomato is an important agricultural product that makes a significant contribution to global food production [2].

Despite the diversity and importance of tomatoes, tomatoes face many disease threats, such as late blight caused by Phytophthora infestans and fungal diseases that cause molds. The impact of these diseases on the tomato crop poses a constant risk to agriculture and can lead to severe economic and food shortages [3] [4]. Timely and accurate identification of these diseases is important for effective control and prevention [5] [6].

The method of detecting tomato diseases has always relied on manual inspection and laborious experiments that have their limitations. Depending on the interpreter, visual inspection may not be necessary for early detection of disease [7] [8] [9] [10] [11] [12]. While the tests are accurate, they are time-consuming and expensive, making it difficult to meet the urgent needs of today's agriculture. To solve these problems, artificial intelligence (AI) has been transformed into agriculture in recent years, especially thanks to advances in deep learning.

Deep learning is a category of machine learning that uses multiple layers of artificial neural networks to solve complex problems. In deep learning, convolutional neural networks (CNN) have become powerful tools for data visualization and have become relevant in many fields [13] [14] [15] [16]. CNNs are characterized using convolutional techniques and are good at extracting relevant features from images, making them ideal for tasks such as image recognition and classification.

This change also extends to agriculture, where CNNs provide a way to quickly and accurately identify diseases, pests, and overall crop health. While this change applies to many crops, we focus on early detection of tomato leaf diseases. Against this problem, ResNet-18 architecture stands out as a light source that measures computational power and accuracy. Our research uses the ResNet-18 architecture to explore the potential of deep learning to transform the tomato plants system and even the permaculture system [17] [18] [19] [20] [21].

Our work set out on a journey to combine the fundamental world of tomatoes with the transformative power of deep learning. The stakes are high as we try to provide fast, effective solutions to ongoing challenges like growing healthy tomatoes, benefiting farmers, fields, and people around the world who depend on this versatile and important fruit. In the following sections, we will describe the process, results and conclusions of our research, which leads to a general discussion about permaculture and its important role in shaping the future of intelligence [22][23][24][25].

## II. LITERATURE SURVEY

These days, there is a lot of research being done in the broad field of image processing applications for plant disease detection and classification. For the prompt identification of plant diseases, these applications are helpful. For any plant, diseases like fungus, bacteria, and viruses can be fatal. In Sabrol et al.'s [26] study, tomato late blight, Septoria spot, bacterial spot, bacterial canker, tomato leaf curl, and healthy tomato plant leaf and stem images are categorized into five categories. The categorization process involves removing

characteristics related to color, shape, and texture from images of healthy and unhealthy tomato plants. Following the segmentation step comes the feature extraction. It is extracted and put into the classification tree from segmented images.

A three-compact convolutional neural network (CNN) pipeline for the automatic detection of tomato leaf diseases has been proposed by Attallah et al. [28]. In order to obtain a more streamlined and sophisticated representation, deep features are extracted from the last fully connected layer of the CNNs using transfer learning. Then, in order to take use of each CNN structure, it combines features from the three CNNs. It then selects and creates an extensive feature set of smaller dimensions using a hybrid feature selection technique. The process for identifying tomato leaf diseases uses six classifiers. In order to confirm the suggested pipeline's ability to compete, the experimental results are also compared with earlier studies on the classification of tomato leaf diseases.

The circumstances of a tomato plant have been determined using a basic CNN model that contains eight hidden layers. When compared to other traditional models, the suggested strategies produce optimal results [24] [27] [29] [30] [31]. The image processing system recognizes and categorizes tomato plant illnesses using deep learning techniques. Here, the author implemented a full system using CNN and the segmentation technique. To attain higher accuracy, a CNN model modification has been implemented.

TABLE I.  TABLE SUMMARIZING THE LIMITATIONS AND GAPS IN CURRENT RESEARCH ON IMAGE PROCESSING APPLICATIONS FOR PLANT DISEASE DETECTION AND CLASSIFICATION

| Limitation | Description |
|---|---|
| Limited Dataset Diversity | o  Focus on specific diseases or crops. <br> o  Lack of diversity may hinder model robustness. |
| Generalization to New Diseases | o  Models might struggle with novel, unseen diseases. |
| Transferability across Crops | o  Investigating the generalizability is essential. <br> o  Models designed for one crop may not adapt well to others. <br> o  Assessing cross-crop transferability is crucial |
| Robustness to Environmental Factors | o  Impact of environmental variations on model performance. <br> o  Models may need to handle real-world, uncontrolled conditions. |
| Interpretability of Deep Learning Models | o  Complex models like CNNs often lack interpretability. <br> o  Understanding model decisions is important for user trust. |
| Data Annotation Challenges | o  Manual annotation challenges affect dataset quality. <br> o  Exploring improved annotation methods is necessary. |
| Real-time Application Challenges | o  Computational efficiency and hardware constraints for real-time deployment. <br> o  Exploring lightweight architectures and edge computing solutions. |

LeNet has been applied to the identification and classification of tomato illnesses while requiring the least amount of CPU processing power. Moreover, to increase classification accuracy, the automatic feature extraction

technique has been used [32]. Table I provides a comprehensive summary of the limitations and gaps identified in current research on image processing applications for plant disease detection and classification.

The work in this paper supports the ResNet-18 architecture and data preparation techniques on a well-organized dataset. With an impressive 91% F1 score, the model demonstrates its ability to differentiate between good and unhealthy tomato leaves. The model's outstanding performance is demonstrated by metrics like accuracy, sensitivity, specificity, and a high AUC score on the ROC curve. This work is significant because it has real-world applications for early disease identification in agriculture. The ResNet-18 model is a potent instrument for crop management that contributes to sustainable agriculture and global food security because of its high precision and specificity.

## III. METHODS

### A. Datasets

The foundation of our classification model lies in the dataset of tomato leaf images, which we've organized into two primary categories: 'healthy' and 'unhealthy.' These categories represent the key labels for our classification task.

*1) Healthy data split:* Within the 'healthy' category, we further partitioned the dataset into 'train' and 'test' subsets. The 'healthy train' subset comprises a total of 1491 images. This subset is instrumental in training the model to recognize healthy tomato leaves effectively. It forms the basis for teaching the model the visual characteristics of undisturbed, healthy foliage. The 'healthy test' subset, on the other hand, consists of 100 distinct images. This collection serves as an independent assessment tool to gauge the model's performance. These images, being separate from the training data, help us evaluate how well the model generalizes to unseen healthy leaf samples.

*2) Unhealthy data split:* The 'unhealthy' category is a bit more nuanced. It encompasses tomato leaves affected by various forms of blight, encompassing both early and late stages of the disease. These stages are categorized as 'unhealthy' for our classification purpose. The 'unhealthy train' subset contains 2809 images. This vast and diverse collection offers a comprehensive training experience for the model to learn the intricacies of identifying blight in its various forms. The inclusion of both early and late stages ensures that the model grasps the entire spectrum of blight-related visual cues. The 'unhealthy test' subset, consisting of 100 images, is divided further into 29 images representing the early stages of blight and 71 images depicting the late stages. This differentiation helps us evaluate the model's proficiency in distinguishing between the progression of blight, which is a critical aspect of disease classification in the agricultural context.

*3) Data validation during training:* It's imperative to ensure that our model generalizes well and doesn't over fit to the training data. To address this concern, we implement a

validation procedure during training. Specifically, we reserve 15% of the data in each training batch for validation. This data is selected randomly in a way that ensures its independence from the training set. The model is periodically evaluated on this validation data to monitor its performance and make adjustments to the training process as necessary.

### B. Data Preprocessing

Tomato leaf images, like many real-world images, exhibit natural variations due to factors such as lighting, background, and color. To enhance the model's ability to classify tomato leaves accurately, we employ several preprocessing steps.

One of the critical preprocessing steps is the application of a mild blurring filter. This filter is designed to reduce noise in the images, resulting in a cleaner and more consistent dataset. By doing so, we aim to minimize the impact of minor variations in image quality that may not be indicative of the actual health state of the tomato leaf. Fig. 1 illustrates the result of applying the blurring filter with a 5x5 kernel, highlighting the smoothing effect on the image.

In addition to blurring, we employ data augmentation techniques to augment the training dataset. Data augmentation is a strategy that introduces variability into the training data by applying random transformations to the images. This process helps the model generalize better, as it exposes the model to a wider range of scenarios and variations during training.

The data augmentation techniques used includes random rotations and flips. Random rotations introduce diversity by rotating images by various degrees, simulating the variability in the orientation of tomato leaves in real-world scenarios. Flips, both horizontally and vertically, add further diversity by reflecting images, mirroring the possible orientations of leaves. This augmentation is shown in Fig. 2. These preprocessing and data augmentation steps collectively serve to enhance the robustness of our classification model. They enable the model to better handle the inherent variations in real-world images, leading to improved accuracy and generalization in the task of tomato leaf classification.



Fig. 1.    Result of blurring filter with a 5x5 kernel.



Fig. 2.    Result of random rotations and flips.

## IV. MODEL ARCHITECTURE

In this research, we opt for the ResNet-18 architecture as the backbone for our classification model. ResNet, short for Residual Network, is a class of neural network architectures that has demonstrated remarkable effectiveness in various deep learning tasks, particularly in image classification. What sets ResNet apart is its ability to address the vanishing gradient problem, a common challenge when training deep networks, through the use of skip connections.

ResNet-18 is a specific variant of the ResNet architecture that we have chosen for our task. It strikes a well-balanced compromise between model depth and computational efficiency. This balance is crucial, especially in applications where both high accuracy and efficient processing are essential. ResNet-18's architecture is structured in a way that enables it to capture intricate features from images, making it an ideal choice for our task of tomato leaf classification. Fig. 3 depicts the Resnet 18 architecture.

ResNet-18 constitutes a convolutional neural network (CNN) architecture specifically tailored for image classification tasks. Noteworthy for its ability to accommodate input images of 224x224 pixels, this architecture serves as an exemplar in the ResNet lineage, embodying the principles of residual learning to facilitate the training of deep networks.

### A. Convolutional Layers

The initial layer of ResNet-18 is a conventional convolutional layer, employing a 7x7 kernel and 64 filters. This layer is succeeded by batch normalization and rectified linear unit (ReLU) activation. Subsequently, a pivotal spatial reduction is accomplished through a 3x3 max-pooling layer with a stride of 2, optimizing the network's capacity to process spatial features.

### B. Residual Blocks

ResNet-18 features four residual blocks, each composed of two consecutive convolutional layers. Within these blocks, each convolutional layer is succeeded by batch normalization and ReLU activation. The hallmark of ResNet architecture, the inclusion of residual connections, enables the creation of shortcut paths, allowing information to bypass one or more layers. This strategic integration mitigates the vanishing gradient problem, empowering the network to learn more effectively.

### C. Global Average Pooling (GAP)

Following the residual blocks, a Global Average Pooling (GAP) layer is applied. This layer computes the average value of each feature map, ensuring a fixed-size output independent of the input dimensions. This pooling operation contributes to the model's spatial invariance and parameter reduction, paving the way for more efficient processing in subsequent layers.

### D. Model Training

To prepare our model for the specific task of classifying tomato leaves as healthy or unhealthy, we undertake a comprehensive training process. Our model is trained for a total of 100 epochs. The choice of training duration is influenced by the complexity of the problem and the size of our dataset.



Fig. 3. Proposed model architecture with ResNet-18 classifier.

*E. Objective Function Binary Cross-Entropy Loss*

In this binary classification problem, we employ Binary Cross-Entropy (BCE) loss as our objective function. BCE loss is a standard choice for binary classification tasks and is well-suited for distinguishing between healthy and unhealthy tomato leaves. It quantifies the dissimilarity between the predicted class probabilities and the ground truth labels. The equation for BCE Loss is shown in Eq. (1) where N represents the total number of data points or samples in your dataset. In the context of the tomato leaf classification, each image of a tomato leaf (whether healthy or unhealthy) is considered a data point. $y_i$ is the ground truth label for the i-th data point. In binary classification, this is a binary value indicating the actual class of the data point. For instance, $y_i = 1$ might represent an unhealthy leaf, and $y_i = 0$ might represent a healthy leaf. $\hat{y}_i$ is the predicted output or probability assigned by the model for the i-th data point. In the context of binary classification, this value represents the model's estimate of the probability that the given data point belongs to the positive class (unhealthy) or negative class (healthy).

$$BCE = -\frac{1}{N} \sum_{i=0}^{N} y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i) \quad (1)$$

We utilize the AdamW optimizer during training, setting the initial learning rate at 1e-4. AdamW is a variant of the Adam optimizer that introduces weight decay, contributing to improved training stability. To further optimize training, we incorporate the Cosine Annealing learning rate scheduler with a Tmax (maximum number of iterations) of 10. This scheduler cyclically adjusts the learning rate, allowing the model to explore different regions of the loss landscape. While a warmup learning rate scheduler was initially implemented during the early training iterations, it was eventually discarded due to providing minimal benefit.

*F. Early Stopping*

To prevent over fitting and to make the training process more efficient, we implement early stopping. This technique introduces a 'patience' parameter, set at 5 epochs in our case. The validation loss serves as the key metric for early stopping. If the validation loss does not improve for five consecutive epochs, the training process is halted.

*G. Model Selection*

Model selection is a critical step in our training process. We choose the top five models with the lowest validation loss and assess their performance. Ultimately, the model with the highest validation F1-score among these top five models is selected as our final model. The F1-score is particularly valuable in binary classification tasks as it balances precision and recall.

*H. Model Testing*

The chosen model is subjected to a comprehensive evaluation process using a range of metrics. These metrics include the F1-score, accuracy, sensitivity, and specificity. The use of multiple metrics allows us to assess the model's performance from different angles and provides a more comprehensive view of its capabilities.

*I. ROC Curve and Operating Point*

To determine the model's operating point, we leverage the Receiver Operating Characteristic (ROC) curve. The ROC curve illustrates the trade-off between sensitivity and specificity at various threshold settings. By analyzing this curve, we can select an operating point that best suits the specific needs of our application. This operating point defines the decision boundary for classifying tomato leaves as healthy or unhealthy. We fine-tune the pre-trained ResNet-18 model on our tomato leaf dataset, using Binary Cross-Entropy (BCE) loss as the objective function. BCE loss is a common choice for binary classification problems, such as distinguishing between healthy and unhealthy tomato leaves.

## V. RESULTS

In our study, we conducted an extensive evaluation of our tomato leaf classification model, and the results demonstrate its exceptional performance in this critical task. The performance metrics we have achieved are truly remarkable and bode well for the practical application of the model. First and foremost, our model exhibited an F1 score of 91%, which is a noteworthy composite metric combining both precision and recall. This F1 score reflects the model's robust ability to accurately classify both healthy and unhealthy tomato leaves. Moreover, our model's accuracy reached an impressive 97%, indicating its proficiency in correctly categorizing tomato leaves.

Precision, which measures the ratio of true positive predictions to the total positive predictions, is a vital metric in binary classification problems. In our case, the model achieved a precision of 90.2%, highlighting its ability to make accurate positive predictions. Additionally, the model showed an outstanding sensitivity of 92%, underscoring its capability to correctly identify unhealthy tomato leaves. Notably, the specificity of the model was also high, at 90%. This indicates that the model was effective in correctly identifying healthy tomato leaves, reducing the risk of false alarms in disease detection systems.

The performance metrics depicted in Table II and Fig. 4 summarizes the evaluation of the two-class classification model for tomato leaf health. The F1 score, a balanced measure of precision and recall, stands at an impressive 91%, indicating the model's robustness. With 97% accuracy, 90.2% precision, 96% sensitivity, and 98% specificity, the model demonstrates high proficiency in distinguishing between healthy and unhealthy tomato leaves, showcasing its reliability in practical agricultural applications.

TABLE II. PERFORMANCE METRICS FOR THE TWO-CLASS CLASSIFICATION PROBLEM

| PERFORMANCE METRIC | VALUE |
|---|---|
| F1 score | 91% |
| Accuracy | 97% |
| Precision | 90.2% |
| Sensitivity | 96% |
| Specificity | 98% |

```
Accuracy: 0.97
Sensitivity: 0.96
Specificity: 0.98
```

Fig. 4.    Metrics for the two-class classification problem.

The visual representation of our model's performance in the confusion matrix (see Fig. 5) offers a more detailed breakdown of its classification accuracy, illustrating the number of true positives, true negatives, false positives, and false negatives. Furthermore, the Receiver Operating Characteristic (ROC) curve, depicted in Fig. 6, plays a crucial role in binary classification tasks. It assesses the trade-off between the true positive rate and the false positive rate at various classification thresholds. Our model's substantial Area Under the Curve (AUC) score of 0.92 on the ROC curve is a testament to its ability to effectively balance the need for disease detection while minimizing the risk of false alarms.

These results underscore the significant potential of our model for practical applications in agriculture. The high F1 score, accuracy, sensitivity, specificity, and AUC, coupled with the model's capacity to balance precision and recall, make it a powerful tool for early disease detection and crop health management. In conclusion, our research findings indicate that our classification model is highly effective and has promising implications for the agriculture industry. It provides a valuable and reliable tool for the early identification of tomato leaf diseases, which can profoundly impact crop yield, food security, and the overall health of the global food supply chain.

The figures, presented as Fig. 7 and Fig. 8, visually distinguish a healthy leaf (see Fig. 7) from an unhealthy one (Fig 8). These depictions serve to provide a quick and clear reference, enabling visual recognition of key characteristics associated with leaf health and distress.



Fig. 6.    ROC curve with AUC of 0.92.



Fig. 7.    Healthy leaf.



Fig. 5.    Confusion Matrix for healthy and unhealthy class.



Fig. 8.    Unhealthy leaf.

## VI. Discussion

The remarkable performance of our ResNet-18-based classifier in the task of tomato leaf classification opens up significant avenues for real-world applications, particularly in the field of agriculture. The model's capacity to distinguish between healthy and unhealthy tomato leaves with over 90% accuracy, sensitivity, and specificity is not only a testament to its efficacy but also holds substantial promise for addressing real-world agricultural challenges, especially in the context of early disease detection and mitigation. One of the most noteworthy implications of our results is the potential for early disease detection. The ability to accurately identify unhealthy tomato leaves with a high degree of sensitivity allows farmers and agricultural professionals to take swift and targeted actions. Early detection of diseases such as blight can lead to more efficient intervention strategies, reducing crop losses, and minimizing the need for extensive pesticide application. This not only has financial benefits for farmers but also contributes to more sustainable farming practices by reducing the environmental impact of pesticide usage.

The high precision of our model, with a precision score of 90.2%, is a critical aspect of its practical utility. It implies that when our model makes a positive prediction (i.e., classifies a leaf as unhealthy), it is overwhelmingly likely to be accurate. This reliability is paramount when it comes to making decisions about crop management and implementing disease control measures. Farmers can have confidence in the model's assessments and act promptly to protect their crops. Equally significant is the specificity of our model, which, at 90%, demonstrates its ability to correctly identify healthy tomato leaves. This means that the risk of false alarms, where healthy plants are incorrectly identified as unhealthy, is minimal. Again, this aspect is crucial in practical agricultural applications where false alarms can lead to unnecessary actions and resource allocation. The high Area Under the Curve (AUC) score obtained on the Receiver Operating Characteristic (ROC) curve is of paramount importance. It signifies the model's proficiency in distinguishing between healthy and unhealthy tomato leaves while maintaining a low rate of false positives. In practical terms, this means that our model effectively strikes a balance between the need for disease detection and the avoidance of unnecessary interventions. The substantial AUC score reassures farmers and agricultural professionals that the model's predictions are both accurate and reliable.

## VII. Conclusion

In summary, our research showcases the immense potential of deep learning and convolutional neural networks in addressing pressing challenges in agriculture, with a specific focus on early disease detection in tomato crops. This technology is not confined to tomatoes alone and can be extended to various other crops, offering invaluable insights and support for sustainable agricultural practices. To harness this potential fully, future work should focus on the practical deployment of these models, integrating them with smart agriculture systems that enable timely responses to disease outbreaks, thus ensuring global food security and promoting sustainable agriculture.

Our observations are a reflection of the robustness of our classification model. The high accuracy, sensitivity, specificity, precision, and AUC score are the result of a well-trained model that has learned to recognize the subtle visual cues associated with healthy and unhealthy tomato leaves. The data preprocessing steps, including blurring and data augmentation, have contributed to the model's ability to handle real-world variations in leaf images. Additionally, the choice of the ResNet-18 architecture played a pivotal role in the model's success. ResNet architectures, known for their skip connections, are adept at training deep neural networks effectively. ResNet-18, in particular, struck a balance between model depth and computational efficiency, making it suitable for our classification task. The success of our model underscores the potential of AI and deep learning in addressing agricultural challenges, offering a valuable tool for farmers and researchers to enhance crop management, reduce disease-related losses, and contribute to more sustainable and efficient agricultural practices. The balance between accuracy, sensitivity, and specificity, along with the AUC score, emphasizes the model's real-world applicability, making it a promising asset for the agriculture industry.

## References

[1] Kovalskaya N., Hammond R.W. Molecular biology of viroid–host interactions and disease control strategies. Plant Sci. 2014;228:48–60. doi: 10.1016/j.plantsci.2014.05.006

[2] Ali H., Lali M.I., Nawaz M.Z., Sharif M., Saleem B.A. Symptom based automated detection of citrus diseases using color histogram and textural descriptors. Comput. Electron. Agric. 2017;138:92–104. doi: 10.1016/j.compag.2017.04.008.[

[3] Wilson C.R. Plant pathogens–the great thieves of vegetable value; Proceedings of the XXIX International Horticultural Congress on Horticulture Sustaining Lives, Livelihoods and Landscapes (IHC2014); Brisbane, Australia. 17–22 August 2014.

[4] Peet M.M., Welles G.W. Greenhouse tomato production. *Crop. Prod. Sci. Hortic.* 2005;**13**:257–304.

[5] Zhang S.W., Shang Y.J., Wang L. Plant disease recognition based on plant leaf image. J. Anim. Plant Sci. 2015;25:42–45.

[6] Ma J., Du K., Zheng F., Zhang L., Sun Z. A segmentation method for processing greenhouse vegetable foliar disease symptom images. Inf. Process. Agric. 2019;6:216–223. doi: 10.1016/j.inpa.2018.08.010.

[7] Basavaiah J., Anthony A.A. Tomato Leaf Disease Classification using Multiple Feature Extraction Techniques. Wirel. Pers. Commun. 2020;115:633–651. doi: 10.1007/s11277-020-07590-x.

[8] Adhikari S., Shrestha B., Baiju B., Kumar S. Tomato plant diseases detection system using image processing; Proceedings of the 1st KEC Conference on Engineering and Technology; Laliitpur, Nepal. 27 September 2018; pp. 81–86.

[9] Agarwal M., Singh A., Arjaria S., Sinha A., Gupta S. ToLeD: Tomato leaf disease detection using convolution neural network. Procedia Comput. Sci. 2020;167:293–301. doi: 10.1016/j.procs.2020.03.225.

[10] Ishak S., Rahiman M.H., Kanafiah S.N., Saad H. Leaf disease classification using artificial neural network. J. Teknol. 2015;77:109–114. doi: 10.11113/jt.v77.6463. M

[11] Khan S., Narvekar M. Novel fusion of color balancing and superpixel based approach for detection of tomato plant diseases in natural complex environment. J. King Saud Univ.-Comput. Inf. Sci. 2020 doi: 10.1016/j.jksuci.2020.09.006.

[12] Sabrol H., Satish K. Tomato plant disease classification in digital images using classification tree; Proceedings of the International Conference on Communication and Signal Processing (ICCSP); Melmaruvathur, India. 6–8 April 2016; pp. 1242–1246.

[13] Sharma P., Hans P., Gupta S.C. Classification of plant leaf diseases using machine learning and image preprocessing tech-niques;

Proceedings of the 10th International Conference on Cloud Computing, Data Science & Engineering (Confluence); Noida, India. 29–31 January 2020.

[14] Rangarajan A.K., Purushothaman R., Ramesh A. Tomato crop disease classification using pre-trained deep learning algorithm. Procedia Comput. Sci. 2018;133:1040–1047. doi: 10.1016/j.procs.2018.07.070.

[15] Sangeetha R., Rani M. Tomato Leaf Disease Prediction Using Transfer Learning; Proceedings of the International Advanced Computing Conference 2020; Panaji, India. 5–6 December 2020.

[16] Hasan M., Tanawala B., Patel K.J. Deep learning precision farming: Tomato leaf disease detection by transfer learning; In Proceeding of the 2nd International Conference on Advanced Computing and Software Engineering (ICACSE); Sultanpur, Inida. 8–9 February 2019.

[17] Coulibaly S., Kamsu-Foguem B., Kamissoko D., Traore D. Deep neural networks with transfer learning in millet crop images. Comput. Ind. 2019;108:115–120. doi: 10.1016/j.compind.2019.02.003.

[18] Jiang D., Li F., Yang Y., Yu S. A tomato leaf diseases classification method based on deep learning; Proceedings of the Chinese Control and Decision Conference (CCDC); Hefei, China. 22–24 August 2020; pp. 1446–1450.

[19] Sabrol H., Kumar S. Fuzzy and neural network-based tomato plant disease classification using natural outdoor images. Indian J. Sci. Technol. 2016;9:1–8. doi: 10.17485/ijst/2016/v9i44/92825.

[20] Mortazi A., Bagci U. Automatically designing CNN architectures for medical image segmentation; Proceedings of the International Workshop on Machine Learning in Medical Imaging; Granada, Spain. 16 September 2018; pp. 98–106.

[21] Salih T.A. Deep Learning Convolution Neural Network to Detect and Classify Tomato Plant Leaf Diseases. Open Access Libr. J. 2020;7:12. doi: 10.4236/oalib.1106296.

[22] Agarwal M., Gupta S.K., Biswas K.K. Development of Efficient CNN model for Tomato crop disease identification. Sustain. Comput. Inform. Syst. 2020;28:100407–100421. doi: 10.1016/j.suscom.2020.100407.

[23] Li G., Liu F., Sharma A., Khalaf O.I., Alotaibi Y., Alsufyani A., Alghamdi S. Research on the natural language recognition method based on cluster analysis using neural network. Math. Probl. Eng. 2021;2021:9982305.

[24] Kaur P., Gautam V. Research patterns and trends in classification of biotic and abiotic stress in plant leaf. Mater. Today Proc. 2021;45:4377–4382. doi: 10.1016/j.matpr.2020.11.198.

[25] Wisesa O., Andriansyah A., Khalaf O.I. Prediction Analysis for Business to Business (B2B) Sales of Telecommunication Services using Machine Learning Techniques. Majlesi J. Electr. Eng. 2020;14:145–153. doi: 10.29252/mjee.14.4.145

[26] H. Sabrol and K. Satish, "Tomato plant disease classification in digital images using classification tree," 2016 International Conference on Communication and Signal Processing (ICCSP), Melmaruvathur, India, 2016, pp. 1242-1246, doi: 10.1109/ICCSP.2016.7754351.

[27] Kaur P., Gautam V. Plant Biotic Disease Identification and Classification Based on Leaf Image: A Review; Proceedings of the 3rd International Conference on Computing Informatics and Networks (ICCIN); Delhi, India. 29–30 July 2021; pp. 597–610.

[28] Suryanarayana G., Chandran K., Khalaf O.I., Alotaibi Y., Alsufyani A., Alghamdi S.A. Accurate Magnetic Resonance Image Super-Resolution Using Deep Networks and Gaussian Filtering in the Stationary Wavelet Domain. IEEE Access. 2021;9:71406–71417.

[29] Wu Y., Xu L., Goodman E.D. Tomato Leaf Disease Identification and Detection Based on Deep Convolutional Neural Net-work. Intelli. Autom. Soft Comput. 2021;28:561–576.

[30] Tm P., Pranathi A., SaiAshritha K., Chittaragi N.B., Koolagudi S.G. Tomato leaf disease detection using convolutional neural networks; Proceedings of the Eleventh International Conference on Contemporary Computing (IC3); Noida, India. 2–4 August 2018; pp. 1–5.

[31] Kaushik M., Prakash P., Ajay R., Veni S. Tomato Leaf Disease Detection using Convolutional Neural Network with Data Augmentation; Proceedings of the 5th International Conference on Communication and Electronics Systems (ICCES); Coimbatore, India. 10–12 June 2020; pp. 1125–1132

[32] Trivedi NK, Gautam V, Anand A, Aljahdali HM, Villar SG, Anand D, Goyal N, Kadry S. Early Detection and Classification of Tomato Leaf Disease Using High-Performance Deep Neural Network. Sensors (Basel). 2021 Nov 30;21(23):7987. doi: 10.3390/s21237987. PMID: 34883991; PMCID: PMC8659659.

# EmotionNet: Dissecting Stress and Anxiety Through EEG-based Deep Learning Approaches

Yassine Daadaa

College of Computer and Information Sciences,
Imam Mohammad Ibn Saud Islamic University (IMSIU), Riyadh, Saudi Arabia

*Abstract*—**Amid global health crises, such as the COVID-19 pandemic, the heightened prevalence of mental health disorders like stress and anxiety has underscored the importance of understanding and predicting human emotions. Introducing "EmotionNet," an advanced system that leverages deep learning and state-of-the-art hardware capabilities to predict emotions, specifically stress and anxiety. Through the analysis of electroencephalography (EEG) signals, EmotionNet is uniquely poised to decode human emotions in real time. To get information from pre-processed EEG signals, the EmotionNet architecture combines convolutional neural networks (CNN) and long short-term memory (LSTM) networks in a way that works well together. This dual approach first decomposes EEG signals into their core alpha, beta, and theta rhythms. We preprocess these decomposed signals and develop a CNN-LSTM-based architecture for feature extraction. The LSTM captures the intricate temporal dynamics of EEG signals, further enhancing understanding. The end process discerningly classifies signals into "stress" or "anxiety" states through the AdaBoost classifier. Evaluation against the esteemed DEEP, SEED, and DASPS datasets showcased EmotionNet's exceptional prowess, achieving a remarkable accuracy of 98.6%, which surpasses even human detection rates. Beyond its technical accomplishments, EmotionNet emphasizes the paramount importance of addressing and safeguarding mental health.**

*Keywords—Electroencephalography (EEG); Long short-term memory (LSTM); Convolutional neural network (CNN); human stress; anxiety detection; deep learning*

## I. INTRODUCTION

Human emotionality, with its intricate facets, has been a subject of rigorous study and interest for decades [1]. Particularly, emotions like stress and anxiety significantly influence behavior, cognition, and overall well-being. The myriad stimuli that individuals face in their daily lives elicit a plethora of emotional reactions, profoundly anchored to brain activity [2]. Occasionally, misinterpretation of these behavioral shifts can lead to potential misdiagnoses [3]. It is crucial to note that while academic literature often uses 'stress' and 'anxiety' interchangeably due to symptom similarities, clear distinctions exist [4]. Stress usually emanates from external stimuli and can manifest as anger, unhappiness, or feelings of overwhelm. Conversely, anxiety is persistent; lingering even after the causative stressor is resolved, and often marked by symptoms like restlessness, nervousness, or unease [5].

There are nuanced categorizations within anxiety itself, notably state and trait anxiety [6]. State anxiety relates to immediate, situational responses, while trait anxiety is an enduring aspect of an individual's personality. Researchers employ distinct methodologies, such as rest state recordings and responsive tests, to measure these anxiety types [7]. With the prevalence of anxiety disorders affecting a significant portion of the global population, understanding them becomes imperative [8]. Statistics, such as those from the USA, show alarming rates of anxiety disorders and related hospitalizations, emphasizing the need for accurate diagnosis and intervention [9]. Moreover, numerous studies have well documented the correlation between anxiety disorders and other medical conditions, such as cardiovascular diseases [10].

Clinical diagnosis of anxiety poses challenges primarily because of the symptomatic overlap with other conditions like depression [11]. Though symptom-based diagnosis assists clinicians, it doesn't provide an objective, quantifiable measure of the underlying causes. In this context, the biomedical community suggests certain chemical biomarkers as promising tools for anxiety assessment [12]. Emerging technologies promise innovations in emotional analysis. "EmotionNet," an advanced system, leverages the power of LSTM networks and CNNs for detailed stress and anxiety detection through EEG signals [13]. As brain state detection advances, researchers view EEG signal analysis as a transformative tool that offers insights into the brain's electrical activities and corresponding emotions [14]. As neural networks improve, they can process these EEG signals, particularly when transformed into spectrograms, to reveal the intricate details necessary for precise stress identification [15]. Finally, while traditional neural networks have made significant strides, there's a pressing need for more nuanced, advanced systems. Emphasizing relevant feature extraction, considering the challenges of datasets and enhancing accuracy are pivotal. Current methodologies, like spectrogram-based and signal processing-based techniques, offer great promise for refining emotional analysis [16].

The proposed study looked into how EEG data parameters (such as electrode selection and frequency bands) affect the classification of anxiety. However, it had some problems, like not being very good at detecting anxiety levels and having a long feature vector length. In contrast, the proposed approach refines this by selecting an optimal subset of EEG features, ensuring better efficiency without compromising the entire EEG data's breadth. This paper introduces EmotionNet, a novel hybrid architecture that significantly advances the field by discerning emotions from EEG signals. There are the following main contributions to this paper:

*1)* Researchers introduced a unique preprocessing methodology that transformed EEG data into azimuthal projection images. By focusing on the alpha, beta, and theta signals, this method provided a fresh perspective on stress detection, enhancing the richness and specificity of the data.

*2)* Researchers developed a pioneering model called "EmotionNet." This hybrid system, combining the strengths of both LSTM and CNN, processes the azimuthal projection images derived from EEG signals. Its robust architecture classifies these images into two distinct classes: stress and non-stress. This innovative integration stands as a hallmark of blending traditional EEG processing with advanced neural network architectures.

*3)* By leveraging the augmented dataset for both training and testing phases, we achieved a significant enhancement in stress and anxiety detection accuracy. We also compared the system's performance with existing state-of-the-art methods. The results underscored the model's superiority and its potential to set new benchmarks in EEG-based stress detection.

In practice, this research has provided transformative contributions to the domain of stress and anxiety detection using EEG signals, setting new standards in preprocessing, model development, and overall system accuracy.

The organization of this paper is as follows: Section II delves into a comprehensive literature review, setting the groundwork for the study. Section III introduces the proposed methodology and elaborates on the specifics of the model. Section IV compares the results derived from the dataset utilized with established state-of-the-art methods for a comparative understanding. The paper culminates in Section V and Section VI, offering discussion and a concise conclusion reflecting on the study's findings.

## II. LITERATURE REVIEW

Emotion recognition, using EEG signals, has been a focal point in Emotion recognition using EEG signals has been a focal point in various studies. In the study [7], a headband equipped with four screen-printed active electrodes was utilized to capture EEG signals. OpenViBE, an open-source software, processed the EEG signals captured using a headband equipped with four screen-printed active electrodes. Additionally, the signals were amplified through an "EEG-SMT" biofeedback board. We employed classification algorithms such as Signal Power (SP), Power Spectral Density (PSD), and Common Spatial Pattern (CSP). Similarly, in a study [8], the MUSE 2 headset recorded neuro-psychological signals from subjects as they viewed standardized movie clips. An LSTM deep learning model processed this data.

Study [9], utilizes the DEAP dataset, recorded EEG signals from 32 volunteers. Feature extraction focused on temporal, regional, and asymmetric dimensions, with a deep learning classifier aiding in emotion categorization. In studies [10] and [16], participants watched music video clips. For EEG feature extraction, researchers employed wavelet transform and approximate entropy, and for emotion classification, they utilized machine learning classifiers such as SVM and Random Forest. The study in [12] took a multimedia approach, combining EEG with galvanic skin responses to recognize emotions.

The potential of convolutional neural networks (CNN) in this domain was highlighted in a study [17], which introduced a randomized CNN model, significantly reducing the need for backpropagation. This approach, on the DEAP dataset, yielded impressive results. Building on this, the study [18] integrated principles from genetic code, achieving up to 92% accuracy on datasets like DEAP and MAHNOB. A study [19] explored stress's health implications, using the EEGnet model to achieve 99.45% accuracy in detecting stress levels in subjects exposed to music experiments.

Advancing further, study in [19] integrated multi-input CNN-LSTM models to analyze fear levels, while study [20] employed CNNs on the UCI-ML EEG dataset to diagnose alcoholism, achieving a 98% accuracy rate. A study [21] merged deep learning models for stress detection, emphasizing their superiority over traditional models. The MODMA dataset was the foundation for the study [22], which utilized CNNs and recorded a commendable 97% accuracy rate. A study [23] delved into the emotional aftermath of COVID-19 among students using an RCN-L system combined with LightGBM techniques, registering around 92.63% accuracy. Lastly, a study [24] simulated mental stress scenarios in a human-machine context, using neural activation features to achieve an 89% accuracy rate. These studies accentuate the versatility and importance of EEG signals in comprehending emotions, with technology playing a pivotal role in this exploration. The following table shows a summary of the related works as well as their outcomes and the accuracy of the studies as mentioned in Table I.

Many previous works have discussed EEG as a convenient brain imaging technique. Different emotions are the key features used in previous works to determine the accuracy of EEG in emotion detection. Most of the previous work provided a satisfactory accuracy rate. The process of acquiring the relevant signals entails the elimination of noise and artifacts through filtration, and the outcome is analyzed using the frequency domain. Lastly, deep learning will be used to perform all methods of extraction and filtration of the EEG signal, as well as provide a frequency domain to the extracted feature. EGG signals are also applicable in emotion recognition since their devices are available in clinics to aid in the diagnosis of symptoms that are used as data for analysis for further medical interventions. Such applications also help in fostering best practices in the curing and publication of critical medical signal data. The gathering of brainwave signals relies on the electrodes standardized by the EEG signals.

TABLE I.    STATE-OF-THE-ART COMPARISON SYSTEMS USING DEEP LEARNING AND EEG SIGNALS

| Cited Reference | Features | Models | Dataset | Accuracy | Limitations |
|---|---|---|---|---|---|
| [7] | SP, PSD, CSP | OpenViBE, EEG-SMT board | - | 92% | It is used to recognize stress only and computational overloaded |
| [8] | Neuro-psychological signals | LSTM (Deep learning) | Standardized movie clips | Negative and positive emotions | Recognize only negative and positive emotions |
| [9] | Temporal, regional, asymmetric | Fully connected, SoftMax | DEAP | 91% | Computational expensive and not generalized solution. |
| [10] | Wavelet transform, approximate entropy | SVM, Random Forest | 40-minute music videos | - | Recognize emotions positive or negative |
| [12] | EEG, GSR | - | 40 music videos | Arousal, valence, like/dislike, dominance, familiarity | Computational expensive and not generalized solution. |
| [13] | SP, PSD, CSP | OpenViBE, EEG-SMT board | - | 94% | Limited dataset. |
| [14] | Wavelet transform, approximate entropy | SVM, Random Forest | 40-minute music videos | 95% | Recognize Stress and Computational expensive. |
| [15] | - | Randomized CNN | DEAP | At least 95% | Backpropagation can be computationally expensive |
| [25] | Brain rhythm code features | Four conventional ML classifiers | DEAP, MAHNOB, SEED | 78%-92% | Complexity of emotion recognition |
| [19] | EEGnet, mother wavelet decomposition | EEGnet (CNN with Relu) | Music experimentation | 99.45% | Just recognized stress and Computational expensive |
| [18] | Multichannel EEG, peripheral physiological | Multi-Input CNN-LSTM | DEAP | 98.79% | Computational expensive and not generalized solution. |
| [26] | EEG signals | CNN | UCI-ML EEG dataset | 98% | Complexity of EEG signals |
| [21] | DWT-based multi-channel EEG | DWT-based CNN, BiLSTM, 2 layers GRU | - | Better than other models | Computational expensive and not generalized solution. |
| [22] | Multiband EEG | CNN | MODMA | 97% | Not mentioned |
| [23] | EEG signals | RCN-L, LightGBM | Post-COVID-19 emotions | 0.9263 (92.63%) | Emotions impacted by COVID-19 |
| [24] | EEG power spectral density (PSD) | Multiple attention-based CNN | Virtual UAV task | 89.49% (arousal), 89.88% (valence) | Computational expensive and not generalized solution. |

## III.    METHODOLOGY

As emotions are the cause of many diseases, identifying these emotions is crucial in order to get the correct medications. One way of identifying these emotions is by using EEG signals. EEG captures scalp electrical activity generated by brain structures [14]. There are many different devices that capture these electrical activities, e.g., brainwaves or TGAM. These devices can then process the captured signals and extract the desired emotion. Therefore, the proposed study tries to study EEG signals and how to use these devices to get these signals. Then, the ensemble-based deep learning architecture is used to predict the mental status of the user from the EEG signals used as data that are gathered from the device. The system architecture of the proposed EmotionNet is shown in Fig. 1.

### A. Data Acquisition

The proposed approach utilizes the SJTU emotion EEG dataset (SEED) from the brain-like computing and machine learning (BCMI) methods [27]. This dataset features EEG data from 15 subjects, recorded over three sessions as they watched various Chinese film clips eliciting distinct emotions. After each clip, participants shared their emotional responses through questionnaires. The EEG data, captured using a 62-channel electrode cap, was down-sampled to 200 Hz and subjected to a

0-75 Hz band-pass filter. We used the DEAP dataset to analyze emotions through EEG signals [28]. This data encompasses 32 participants exposed to 40 one-minute music videos, each inducing a consistent emotion. Recorded data from 32 EEG channels was down-sampled to 128 Hz for reduced system complexity. Another study [29] employed the DASPS database, which centers on EEG responses during exposure therapy, a variant of cognitive behavioral therapy (CBT). This database comprises EEG data from 23 healthy participants. These participants, prior to the experiment, provided written consent and gauged their anxiety using the Hamilton Anxiety Rating Scale.

The aim of the study was to detect stress using the proposed model. The data was provided in the form of 'mat' files, which were read into the Python program using the SciPy library. This paper used data augmentation techniques to generate new data for training the neural network. In this study, the anxiety state from the SEED, DEAP, and DASPS datasets was considered to be a stressful state for the target task. Data-augmentation techniques can be used to increase the size of the existing EEG dataset. Generating additional data by applying transformations to the existing data, such as shifting the signal in time or adding noise, is performed to increase the sample size. This can help increase the variability in the data and improve the generalizability of the model.

Fig. 1. The system architecture of proposed emotionnet.

## B. Preprocessing

Electroencephalography (EEG) is a modern electrophysiological screening mechanism that is used to record the electrical activities of the brain. The EEG method measures fluctuations in voltage ensuing from the current generated by the flow of ions in the brain neurons [14]. The EEG signals can be categorized into five main groups, which are Delta, Alpha, Theta, Beta, and Gamma. A delta signal or wave is a neural oscillation with high amplitude and varied frequencies ranging from 0.5 hertz to 4 hertz [17]. The wave is commonly associated with sleep. Alpha signals have frequencies ranging from 7.5 hertz to 13 hertz. It is commonly experienced in the posterior areas of the brain when a patient closes and relaxes their eyes. The theta signal is a slow-activity wave with a frequency ranging between 3.5 hertz and 7.5 hertz. It is a normal occurrence in children from 0 to 13 years old, but it indicates sub-cortical lesions, hydrocephalus, or metabolic encephalopathy. The brain exhibits a beta signal when it is aroused and actively engages in activities. It has a frequency of 14 to 35 hertz. Gamma signals indicate that an individual has attained peak concentration and help in information processing. It has a frequency of 35 hertz or more.

*1) Signal filtration:* Signals from an EEG device usually have a lot of noise and other artifacts that may originate from sources that can be biological or environmental [30]. A filter

removes some of the unwanted signal features when processing a signal. Filtering represents a class in signal processing that entails partial or complete suppression of certain aspects of a signal. EEG commonly refers to digital filtering as the usual pre-processing phase in analyzing the EEG data. The usual exercise in processing EEG signals includes applying a high-pass filter for the elimination of the slow frequencies with a lesser amount of 0.1 Hz and a low-pass filter to remove frequencies that are above 40 to 50 Hz.

Signal filtering refers to the modification of a measured signal through the use of an algorithm or logic to eliminate its undesirable features before it is adopted by a controller. Some of the examples in control include feedback variables for proportional-integral-derivative (PID) and advanced process control (APC) controllers [31]. The examples of calculations entail computations centered on steady-state material balances, the process, and the control metrics. The primary objective of signal filtering is to reduce and smooth the high-frequency noise related to flow, temperature, or pressure measurements. Noise related to differential pressure across the orifice plate is a common example used to infer flow rate. High-frequency noises are usually considered random and an additive in the measured signal and are normally uncorrelated in the period Fig. 2.



Fig. 2. a) Anxiety without a filter, Fig. 2(b) Anxiety with a filter, Fig. 2(c) Stress without a filter, and Fig. 2(d) Stress with a filter.

Filters can be categorized based on their design as either finite impulse response (FIR) or infinite impulse response (IIR) [32]. The impulsive response refers to how the filter works with the unit impulse signal within the time domain. The FIR filter has a finite-distance impulse response; then, its output drops to zero, producing equal delays for all frequencies. The IIR filters, on the other hand, have an infinite impulse reaction. It also produces unequal delays. However, its main advantage is that it is computationally highly efficient. Another feature in the design of filters is the signal direction when used as an input. Causal filters comprise past and present information. Similarly, it refers to filters that rely on future and past input as noncausal filters.

After recording and filtering an EEG signal, researchers need to extract its features. There are several methods for

extracting features from an EEG signal. During frequency domain analysis, oscillating parts are used to break down EEG signals and separate out specific neuronal activity. When decomposing time-domain signals into weighted cosine and sine functions, the frequency domain is primarily utilized.

---

**Algorithm 1:** EEG Preprocessing

Input: Raw EEG data: Raw EEG data: R= {r₁, r₂,..., rₙ}

Output: Preprocessed EEG signal: P

Variables: N={n₁,n₂,...,nₘ} : Detected noise in the EEG data. F= {f₁, f₂, ..., fₙ}: Data after filtering. E={e₁, e₂,..., eₙ} : Extracted features from the EEG data. K={k₁, k₂,..., kₙ} : Data structured for KNN classification

Functions: L(R) : Loads the raw EEG data. G(Rᵢ) : Filters segment Ri. T(Fᵢ): Applies Fast Fourier Transform to segment Fᵢ. C(Eᵢ) : Prepares feature Eᵢ for KNN classification.

> Initialization: R←L(R)
>
> Signal Filtration : ∀i∈{1,2, ..., n} ←*DetectNoise (Nᵢ)←DetectNoise (Rᵢ)*
> $$F_i \leftarrow G(R_i - N_i)$$
>
> Feature Extraction: ∀i∈{1,2,...,n}: Eᵢ←T(Fᵢ)
>
> Data Classification Preparation: ∀i∈{1,2,...,n}: Kᵢ←C(Eᵢ)
>
> Advanced Filtration: ∀i∈{1,2,...,n}: ←AdvancedFilter(Fᵢ)
>
> Final Preprocessing: P=K

End

---

The Fast Fourier Transform (FFT) is a feature extraction method used for extracting the finer details of emotions such as spectral entropy and spectral centroid. FFT extracts these simple features from the alpha, beta, and alpha to gamma frequencies. Seeing that the theta and delta have a very low-frequency range, The FFT method does not require the lower frequencies due to their lack of sufficient information. After filtering the signal and isolating the relevant signals, they remain unidentified and require classification. The KNN algorithm has a majority voting scheme, which will be used to classify the unidentified signals. The algorithm classifies the new data based on the highest number of votes. The majority vote schemes are used instead of the similarity vote schemes because they are less sensitive to the outlier, which aids the FFT since it is a method for extracting the finer details [33].

$$Spectral\ Entropy = H(x) = \sum_{x \in X} x_i \cdot \log 2\ x_i \quad (1)$$

$$Spectral\ Centroid = \frac{\sum_{K=1}^{N} kF[K]}{\sum_{K=1}^{N} F[K]} \quad (2)$$

where, $F[K]$ is the amplitude corresponding to bin $k$ in the FFT spectrum.

To filter the EEG signal from noise, this paper used the MNE-Python Library. The MNE-Python Library provides algorithms implemented in Python that cover multiple data pre-processing methods to reduce noise from external (environmental) and internal (biological) sources. The two categories of noise filtering strategies are eliminating contaminated data segments and using signal processing techniques to attenuate artifacts. The MNE-Python library provides these two categories at different stages of the pipeline through functions that use automatic or semi-automatic data pre-processing along with interactive plotting capabilities.

The first step of pre-processing entails restricting the signal to a chosen frequency range. The MNE-Python library includes

various filtering algorithms such as low-pass, high-pass, band-stop, and notch filtering. A high-pass filter is used to filter out slow frequencies and high frequencies with a low-pass filter. And the bandpass, where frequencies pass between defined upper and lower frequencies. Band-stop is the inverse of band-pass, where frequencies between upper and lower defined frequencies are rejected. Instances of raw data are filtered using a method that supports both fast fourier transform (FFT) based on finite impulse response (FIR) and finite impulse response (IIR) filters. The standard multiprocessing Python module exposed with the Joblib Python pipeline tool allows for parallel filtering of multiple channels. We will be using the FIR filter in this paper.

The finite impulse response (FIR) [34] filters can have a linear phase, so they have the same delay at all frequencies, while IRR filters cannot. The phase and delay group characteristics are also usually better for FIR filters. FIR filters are much easier to control and are always stable. FIR filters have a well-defined passband, can be converted to minimum-phase, and can be corrected to zero-phase without additional computations. MNE-Python provides FIR filters with 0.16, 0.15, and 0.13 default constant filter delays. Also, it provides two other filters called MNE-C default and minimum-phase. As shown below in Fig. 3, a signal is tested with different types of FIR filters in MNE-Python and a low-pass of 40 Hz. The blue signal represents the original signal without applying any filtration, whereas the orange signal represents the original signal with noise. Other colored signals represent the type of filter used on them, as shown in Fig. 3.



Fig. 3.   MNE-Python FIR filter types.

*C. Features Extraction*

The authors of this paper utilized the CNN-LSTM [35], [36] feature extraction method. CNN proved to be good at extracting signal patterns but had a disadvantage in terms of long-term dependency. LSTM solves the problem by providing an excellent long-term dependency, allowing it to be used as a time series and negating the CNN disadvantage. After being filtered, the signal goes through the CNN-LSTM process, as shown in the figure below. The classification process will receive the signal after it has been filtered and processed

through the CNN-LSTM process. A convolutional neural network (CNN) is a deep learning algorithm with the ability to process images. CNN also proved that it could detect patterns from brainwaves, such as emotions, in a multi-channel EEG recording, which also gives it the ability to process EEG signals. A long-short-term memory (LSTM) is a neural network that can learn based on the predictions of a given problem. A recurrent neural network (RNN) is a network with highly efficient working internal memory for predicting time series. LSTM is just an extension of an RNN cell, which overhauls the disadvantages of RNN.

Since CNN is an image-processing algorithm, this paper is going to change the EEG signal into an image and pass it to CNN. After passing through CNN, it goes through the LSTM for the time series. Combining the CNN and LSTM is essential, as they rely on each other for effective functioning. The LSTM passes the signal to the classifier to identify the emotion (see Fig. 4).



Fig. 4. CNN combined with RNN-LSTM layers.

CNN applied to the pre-processed signals. CNN has three major elements: local sensing fields, weight sharing, and downsampling. These three elements can decrease network complexity, which is good. Also, CNN has high accuracy because it can learn from non-linear convolution and local non-linear activation functions. Many CNNs combine using pooling as layers to create a close enough representation of the intermediate features from the signals, expressing a high level of features. The convolution layer uses a filter on the input data to produce feature maps. The filter slides over the input to execute the convolution. Matrix multiplication is performed at every position, and the results are then summed onto the feature map. CNN's pooling layer takes smaller samples of the features that the convolution layer found. This cuts down on the amount of work that needs to be done and the extent to which the network is overfitted. Only necessary information should be extracted from a pooling process, and irrelevant information should be discarded. This greatly enhances the performance of CNN. Fig. 5 shows the convolution-max pooling process. Fig. 6 shows the proposed CNN model structure.



Fig. 5. Convolution-Max pooling diagram.



Fig. 6. CNN-LSTM structure diagram.

As seen in Fig. 5 and Fig. 6, CNN is made for capturing local spatial features of the data, but CNN cannot capture the data sequence in a long-term dependence relationship, and it can vanquish the weaknesses of CNN. A combination of CNN and LSTM creates a hybrid model, resulting in excellent performance in signal recognition. The raw data is pre-processed and filtered from the noise, and then it enters the CNN model for feature maps before entering the LSTM for the time series.

---

**Algorithm 2:** EEG Feature Extraction using CNN and LSTM

---

Input: Preprocessed EEG signal: Preprocessed EEG signal: $P=\{p_1, p_2,..., p_n\}$

Output: Features: $F=\{f_1, f_2,..., f_n\}$

Variables: $C=\{c_1, c_2,..., c_n\}$ : Features extracted by CNN.
$\quad\quad\quad L=\{l_1, l_2,..., l_n\}$ : Features extracted by LSTM.

Functions: $CNN(P_i)$: CNN model that extracts features from segment $P_i$. $LSTM(P_i)$: LSTM model that extracts features from segment $P_i$.

Combine $(C_i, L_i)$: Combines CNN and LSTM features for segment Pi.

  Initialization: $P \leftarrow$ Load(P)

  Signal Filtration: $\forall i \in \{1, 2,..., n\}$: $\leftarrow$ DetectNoise($N_i$) $\leftarrow$ DetectNoise($R_i$)
  $\quad\quad\quad\quad\quad Fi \leftarrow G(R_i - N_i)$

  Feature Extraction using CNN: $\forall i \in \{1, 2,...,n\}$: $C_i \leftarrow CNN(P_i)$

  Feature Extraction using LSTM: $\forall i \in \{1, 2,...,n\}$: $Li \leftarrow LSTM(P_i)$

  Combining Features: $\forall i \in \{1, 2,...,n\}$: $F_i \leftarrow$ Combine($C_i, L_i$)

End

---

### D. Human Anxiety and Stress Classification (AdaBoost Classifier)

The classification method uses the SoftMax classifier. After the convolution-max pooling has been flattened, it is then passed to the fully connected FC. The SoftMax classifier then receives the final vector as input. The SoftMax classifier then assigns an emotion to the given final vector input. Emotional recognition is the process of identifying emotions. Facial expressions, voice impressions, written texts, psychology, and electrode devices placed on the head can all be used to recognize emotions.

Emotion recognition is going to perform in this paper as follows: A TGAM device is utilized to extract the brain's signals. These signals are called EEG signals and are raw; thus, they need to be cleansed from noise to further increase the accuracy of the emotion extraction. The filtered EEG signal then proceeds to the feature extraction process, where it undergoes a series of methods known as CNN-LSTM. In there, the signal will first go through the convolution and Max-Pooling processes of the CNN several times before being sent to the LSTM for the time series for long-term dependency. Finally, the CNN-LSTM processes the signal and then classifies it using the SoftMax classification method to assign an emotion. After the signal has gone through all these processes, the result will then be that person's emotion at the time of the signal extraction. Because there are many emotions to recognize, this paper focused on detecting anxiety, stress, depression, etc., with good accuracy rather than detecting more emotions with less accuracy.

---

**Algorithm 3:** Human Anxiety and Stress Classification using AdaBoost

---

Input: Extracted Features: F= {$f_1$, $f_2$,..., $f_n$}

Output: Class Labels: L={$l_1$,$l_2$,...,$l_n$} Where $l_i$ can be "Anxiety", "Stress", or "Neutral"

Variables: A= {$a_1$,$a_2$,..., $a_n$}: Classification results using AdaBoost.

Functions: AdaBoost ($F_i$): AdaBoost classifier that determines class label for feature $F_i$.

    Initialization: F← Load(F)

    Anxiety and Stress Classification using AdaBoost: ∀i ∈{1, 2,..., n}: ← $A_i$ ←AdaBoost($F_i$)

    Class Label Assignment: ∀i∈{1,2,...,n}:

    If $A_i$=1 then $L_i$← "Anxiety"

        Else if $A_i$=2 then $L_i$← "Stress"

     Else $L_i$← "Neutral"

End

---

AdaBoost, short for "Adaptive Boosting," has emerged as a potent ensemble machine learning technique that focuses on the principle of amalgamating the strengths of numerous "weak" classifiers to forge a robust classifier. Its application in EEG feature classification for discerning stress and anxiety presents a unique approach that offers notable advantages.

At the beginning of the AdaBoost process, each EEG data point or feature vector weighs equally, ensuring a level playing field. A weak classifier, often a simple decision tree known as a "decision stump," is trained on these features. Despite its designation as "weak," the classifier's aim isn't sheer randomness but to surpass random guesswork, albeit marginally. Following the training, this classifier undergoes an evaluation phase. Meticulously identifying the misclassified instances and incrementing their weights pushes the subsequent classifier to concentrate more assiduously on the challenging, previously misclassified instances. This iterative emphasis on the "hard-to-classify" instances is where AdaBoost truly shines. One of the key steps in the AdaBoost algorithm is the assignment of weights to the classifiers themselves. Classifiers with higher accuracy have greater influence, allowing them to play a more significant role in the final decision-making process. This hierarchy ensures that better-performing classifiers play a pivotal role.

As AdaBoost iterates, the process undergoes fine-tuning. Each cycle refines the classifier weights, focusing more on the problematic instances. Such a continuous feedback loop ensures that, by the end of the specified iterations, the ensemble is adept at handling a majority of the scenarios, including the challenging ones. AdaBoost does not rely on a single classifier to classify a new EEG feature vector. Instead, it consults its ensemble, with each member casting a weighted vote based on its accuracy. The culmination of these votes determines whether the EEG feature vector corresponds to stress, anxiety, or a neutral state. The inclusion of fine-tuning in this process is pivotal. The EEG data, with its intricacies and subtle nuances indicative of stress or anxiety, demands a classifier that is both adaptive and discerning. AdaBoost, with its iterative refinement and emphasis on challenging instances, stands out as an ideal choice. By progressively focusing on the harder-to-classify instances and adjusting classifier influence based on performance, AdaBoost ensures that the final model is not just a mere aggregation but a finely-tuned ensemble primed for accuracy.

## IV. EXPERIMENTAL RESULTS

### A. Experimental Setup

The proposed method is executed on a common platform machine with an Intel Core i7 10th generation processor and 32 GB of RAM without using a GPU. Despite being a complex processing method that combines signal and image processing, it requires a relatively small number of epochs for better accuracy, leading to less training time. It also optimizes the average prediction time for each input. The complete dataset is split into an 80–20 ratio to create the test dataset, with 20% reserved for testing.

### B. Statistical Analysis

Performance evaluation of classification models is vital for understanding their efficacy. In comparing the proposed solutions with existing ones, several metrics are employed.

A fundamental metric for classification models, accuracy provides an aggregate measure of the model's ability to predict correctly. It computes the ratio of correctly predicted instances to the total number of instances, and it's defined as:

$$\text{Accuracy} = (TP+TN)/(TP+TN+FP+FN) \times 100 \quad (3)$$

True Positive (TP) and True Negative (TN) represent correct predictions, while False Positive (FP) and False Negative (FN) denote incorrect predictions. A model with high

accuracy implies reduced prediction errors, which can have significant cost implications.

Often referred to as Recall, True Positive Rate, or Hit Rate, sensitivity captures the model's proficiency in predicting positive instances accurately. It is expressed as:

$$Recall = TP/(TP+FN) \times 100 \qquad (4)$$

A high recall indicates a low FN rate, signifying that the model has a commendable ability to detect positive instances.

This metric evaluates the model's skill in correctly predicting negative instances. Defined as the ratio of true negatives to the sum of true negatives and false positives, it is formulated as:

$$Specificity = \frac{TN}{TN+FP} \times 100 \qquad (5)$$

Fourth measurement is Precision. Precision measures the proportion of accurately predicted positive instances against all predicted positive instances. It's an indication of the model's ability to correctly identify positive instances among all predicted positives. Its formula is:

$$Precision = \frac{TP}{TP+FP} \times 100 \qquad (6)$$

Fifth measurement is F1-score. Combining both precision and recall, the F1-score offers a balanced metric that considers the harmonic mean of precision and recall, making it essential for understanding a classifier's robustness and accuracy. It's defined as:

$$F1 - Score = 2\frac{precision \times recall}{precision + recall} \qquad (7)$$

In essence, these metrics collectively provide a comprehensive view of a classifier's performance, ensuring that its strengths and weaknesses across different dimensions are adequately captured and understood.

*C. Experimental Results*

Python is chosen as the programming language for this paper because it offers simplicity, consistency, flexibility, and accessibility to various libraries and frameworks. Python is a dynamic, high-level object-oriented programming language that offers perfect solutions to machine learning due to its independence. Furthermore, its independence across platforms makes Python resource- and time-saving in deep learning, where the developers would incur more resources to complete a paper. The Python language is reliable due to its ability to run on multiple platforms without the need to change. Python is easy to execute, making it a standalone solution to meet machine learning needs. These features have made it more popular. It's also popular because of its useful libraries and packages that save time and reduce the likelihood of errors. The libraries and frameworks offer a reliable environment for machine learning due to their pre-written codes that speed up coding when working on a complex paper like the current one. Python's interpreted nature allows for faster code execution without the need for a compiler. The aforementioned properties make Python a priority for this paper.



(a)



(b)

Fig. 7. Loss and accuracy of the proposed EmotionNet model, where Fig.7 (a) shows the loss curve and, (b) shows the accuracy curve with 300 epochs.

Fig. 7 appears to present a detailed analysis of the performance metrics for the EmotionNet model for 300 epochs. Fig. 7(a) likely to illustrates the loss curve, which is a graphical representation of how the model's prediction error decreases over time as it learns from the training data. This curve is crucial to understanding how effectively the model is learning and optimizing its parameters. A typical loss curve would show a downward trend, indicating that the model is becoming more accurate in its predictions. Fig. 7(b) probably depicts the accuracy curve, showcasing how the model's prediction accuracy improves across the epochs. The model's proficiency in correctly classifying or predicting emotional states is expected to increase, resulting in an upward trend in this curve. Both of these curves together provide a comprehensive view of the model's learning dynamics and performance, with the loss curve focusing on error minimization and the accuracy curve emphasizing successful predictions.

EmotionNet has obtained accuracy equal to 98.6%. To calculate the metrics of this paper, this study used accuracy, sensitivity (SE), specificity (SP), precision, recall, and F1-Score. Four variables are used in the calculations. These variables are: true positive (TP), which equals 317; true negative (TN), which equals 312; false positive (FP), which equals 4; and false negative (FN), which equals 0. The authors created the confusion matrix below using the 80-20 split for training and testing, respectively, as shown in Fig. 8. In this paper, the author has also calculated the confusion matrix for 70–30, 60–40, and 50–50. As seen from their respective figures, the numbers are much lower than expected. But the accuracy rating is also lower than 80–20. And for that reason, this paper has gone with the 80-20.

(a) 80-20% split.



(b) 70-30 Confusion Matrix.



(c) 60-40 Confusion Matrix.



(d) 50%-50% split.

Fig. 8.   Confusion matrix of different training and testing split ratios.



Fig. 9.   The confusion matrix of detection of stress and anxiety, where 0 shows the detection rate of stress and 1 represents the anxiety by proposed model EmotionNet.

This study has chosen to go with the 80-20 method because it gave us the best result in terms of accuracy. Also, the loss is approximately 11%, which is the least we got. Below are two emotions classes anxiety and stress as shown in Fig. 9. They are both visualized in signals with all eight ranges which are Delta, theta, low alpha, high alpha, low beta, high beta, low gamma, and high gamma. Also, both are shown with and without a filter to showcase the difference between a clean signal and a noisy signal.

As you can see from TABLE II. , the proposed accuracy result is 98.6%. Reference [23] only uses CNN with an accuracy result of 94.83%. Reference [9] only uses LSTM with an accuracy result of 91.85%. This paper used the combination of both CNN and LSTM. It showed that it has better potential rather than just using CNN or LSTM individually. Reference [18] on the other hand uses both CNN and LSTM but has a lower accuracy rating than EmotionNet. This means that having an efficient architecture is most critical. It can be seen from their accuracy result, which is 80.57%.

TABLE II.    STATE-OF-THE-ART COMPARISON

| Ref. | ACC | SE | SP | F1-Score |
|---|---|---|---|---|
| Proposed EmotionNet | 98.6% | 100% | 98.73% | 99.22% |
| LSTM [9] | 91.85% | 94.00% | 96.74% | 95.00% |
| CNN [23] | 94.83% | 86.67% | 98.17% | 89.93% |
| CNN-LSTM [18] | 80.57% | 100% | 71.72% | 76.30% |

## V.    DISCUSSIONS

The exploration and classification of EEG signals to discern and quantify emotional states such as stress and anxiety have witnessed a radical evolution with the integration of advanced machine learning algorithms. At the heart of this investigation is the objective to achieve a nuanced understanding of the myriad emotional responses of the human brain and harness this knowledge for clinical and therapeutic applications.

The initial foray into EEG-based emotion classification was governed by a preliminary preprocessing phase [37]. The preprocessing and filtration stages were crucial in addressing the contamination of EEG recordings by a variety of artifacts, from biological to environmental origins. The defined algorithm effectively trimmed the EEG signal to a desirable frequency range, addressing both high and low frequencies, ensuring an optimized dataset for feature extraction. The adopted approach rigorously eliminated unnecessary complexities and preserved relevant data, laying the groundwork for the subsequent steps.

Table II presents a state-of-the-art comparison of various approaches in the field of emotion recognition. The proposed EmotionNet achieves an impressive accuracy result of 98.6%, showcasing its superiority over other methods. Reference [23] solely employs CNN architecture and achieves an accuracy of 94.83%, while Reference [9] utilizes LSTM and achieves an accuracy of 91.85%. Notably, the proposed EmotionNet combines both CNN and LSTM, demonstrating better potential compared to using CNN or LSTM individually. It is worth mentioning that Reference [18] also employs a combination of CNN and LSTM but achieves a lower accuracy rating than EmotionNet, emphasizing the importance of an efficient architecture. The table includes additional performance metrics such as sensitivity (SE), specificity (SP), and F1-Score for each approach, providing a comprehensive overview of their capabilities in emotion recognition.

Upon having a refined EEG dataset, the challenge transitioned to extracting meaningful features that encapsulate the emotional spectrum of the human brain. This is where the integration of deep learning models, namely CNN and LSTM, came into play. CNNs, with their prowess in handling image-based data, converted EEG signals into spectrogram-based images, enabling a richer feature extraction process. On the other hand, LSTMs processed the sequential data in the time-series nature of EEG data. The symbiosis of CNN and LSTM exhibited efficacy in gleaning relevant features indicative of different emotional states.

However, the pinnacle of exploration was the application of the AdaBoost classifier, fine-tuned to achieve optimal classification results. AdaBoost's adaptability in combining multiple "weak" classifiers to curate a robust classifier became pivotal. Its iterative feedback loop, emphasizing harder-to-classify instances and adjusting weights to improve classification accuracy, offered an adept approach to classifying EEG signals into stress, anxiety, or neutral states. The continuous refinement and fine-tuning of AdaBoost underscored its superiority in handling the intricacies of EEG data.

In summary, the journey from raw EEG data to a nuanced understanding of emotional states has been both intricate and enlightening. Combining preprocessing methods, advanced deep learning models, and adaptive classifiers like AdaBoost showed how EEG data could be used in medical and therapeutic research. As the domain of EEG-based emotion classification expands, the techniques and algorithms outlined in this investigation will inevitably serve as foundational pillars for future research and applications.

The current evaluation utilizes DEEP, SEED, and DASPS datasets. In future iterations, we will train and test EmotionNet on a broader array of datasets to ensure its universality across different demographic and cultural backgrounds. There is potential to integrate EmotionNet into real-time monitoring systems, such as wearable technology, to provide constant mental health feedback and alert individuals or healthcare providers to deteriorating emotional states. While the current accuracy of EmotionNet is commendable, there is always scope for enhancement. Future endeavors can look into refining model parameters, exploring other architectures, or incorporating transfer learning for improved accuracy. EmotionNet's architecture could be adapted to predict a broader spectrum of emotions, expanding its utility in diverse applications, while still maintaining the current focus on stress and anxiety.

In summary, while EmotionNet stands as a significant stride in EEG-based emotion recognition, the journey forward promises further innovation, refinement, and meaningful societal impacts.

## VI. CONCLUSION

The presented work introduces "EmotionNet," a novel deep learning system adept at predicting stress and anxiety levels through EEG signal analysis. The integration of convolutional neural networks (CNN) and long-short-term memory (LSTM) networks serves as a significant advancement in EEG-based emotion recognition. The fact that EmotionNet can achieve a classification accuracy of 98.6% shows how well it works. This is possible by using signal decomposition, preprocessing, and the CNN-LSTM architecture for feature extraction. Furthermore, evaluation of well-regarded datasets like DEEP, SEED, and DASPS reinforces its robustness and reliability in predicting emotional states. EmotionNet not only epitomizes technical progression in the domain but also underscores the broader societal imperative of understanding and prioritizing mental health, especially in times of global challenges like the COVID-19 pandemic.

## REFERENCES

[1] D. Alamsyah and P. Merdeka, "Mental Health as Common Lifestyle.," Journal of Litterature Language and Academic Studies, vol. 2, no. 2, pp. 51-56, 2023.

[2] C.-Y. Liao, . R.-C. Chen and S.-K. Tai, "Emotion stress detection using EEG signal and deep learning technologies.," in IEEE International Conference on Applied System Innovation, Taiwan, 2018.

[3] J. DeMartini, G. Patel and T. L. Fancher, " Generalized anxiety disorder," Annals of internal medicine, vol. 170, no. 7, pp. 49-64, 2019.

[4] S. Mane and A. Shinde, "StressNet: Hybrid model of LSTM and CNN for stress detection from electroencephalogram signal (EEG)," Results in Control and Optimization, vol. 11, p. 100231, 2023.

[5] B. Roy, L. Malviya, R. Kumar, S. Mal and A. Kumar, "Hybrid Deep Learning Approach for Stress Detection Using Decomposed EEG Signals," Diagnostics, vol. 13, no. 11, p. 1936, 2023.

[6] A. Iyer, S. Das, R. Teotia, S. Maheshwari and R. Sharma, "CNN and LSTM based ensemble learning for human emotion recognition using EEG recordings," Multimedia Tools and Applications, vol. 82, no. 4, pp. 4883-4896, 2023.

[7] N. Pusarla, A. Singh and S. Tripathi, ". Normal inverse Gaussian features for EEG-based automatic emotion recognition," IEEE Transactions on Instrumentation and Measurement, vol. 71, pp. 1-11, 2022.

[8] A. Sakalle, P. Tomar, H. Bhardwaj, D. Acharya and A. Bhardwaj, "A LSTM based deep learning network for recognizing emotions using wireless brainwave driven system.," Exoert Systems with Applications, vol. 173, pp. 1-17, 2021.

[9] C. Heng, L. Aiping, Z. Xu, C. Xiang, W. Kongqiao and C. Xun, "EEG-based emotion recognition using an end-to-end regional-asymmetric convolutional neural network," Knowledge-Based Systems, vol. 205, pp. 1-7, 2020.

[10] Z. Jianhua, Y. Zhong, C. Peng and N. Stefano, "Emotion recognition using multi-modal data and machine learning techniques: A tutorial and review," Information Fusion, vol. 59, pp. 103-126, 2020.

[11] Y. Cimtay, E. Ekmekcioglu and S. Caglar-Ozhan, "Cross-Subject Multimodal Emotion Recognition Based on Hybrid Fusion," IEEE Access, vol. 8, pp. 168865-168878, 2020.

[12] S. Koelstra, C. Muhl, M. Soleymani, Y. Lee, Jong-Seok, E. Ashkan, P. Touradj, N. Thierry, P. Anton and P. Ioannis, "Sander, Koelstra., Christ1-ian, M¨uhl., Mohammad, Soleymani., Jong-Seok, Lee., Ashkan, Yazdani., Touradj, Ebrahimi., et al. DEAP: A Database for Emotion Analysis using Physiological Signals," IEEE Transaction on Affective Computing, vol. 3, no. 1, pp. 18-31, 2011.

[13] Y. Wei, Y. Wu and J. Tudor, "A real-time wearable emotion detection headband based on EEG measurement," Sensors and Actuators A Physical, vol. 263, pp. 614-621, 2017.

[14] T. Michael, "Fundamentals Of EEG Measurement," Measurement Science Review, vol. 2, pp. 1-9, 2002.

[15] S. Ibrahim, R. Djemal and A. Alsuwailem, "Electroencephalography (EEG) signal processing for epilepsy and autism spectrum disorder diagnosis," Biocybernetics and Biomedical Engineering, vol. 38, no. 1, pp. 16-26, 2018.

[16] Z. Yaqing, C. Jinling, T. Jen Hong, C. Yuxuan, C. Yunyi, L. Dihan, Y. Lei, S. Jian, H. Xin and C. Wenliang, "An Investigation of Deep Learning Models for EEG-Based Emotion Recognition. Front. Neurosci. 14:622759.," Frontiers in Neuroscience, vol. 14, 2020.

[17] D. Jude Hemanth, "EEG signal based Modified Kohonen Neural Networks for Classification of Human Mental Emotions," Journal of Artificial Intelligence and Systems, vol. 2, pp. 1-13, 2020.

[18] N. Masuda and I. Yairi, "Multi-Input CNN-LSTM deep learning model for fear level classification based on EEG and peripheral physiological signals," Frontiers in Psychology, vol. 14, p. 1141801, 2023.

[19] S. Bhatnagar, S. Khandelwal, S. Jain and H. Vyawahare, "A deep learning approach for assessing stress levels in patients using electroencephalogram signals," Decision Analytics Journal, vol. 7, p. 100211, 2023.

[20] C. Vartak and L. Jolly, "Alcoholic Addiction Detection Based on EEG Signals Using a Deep Convolutional Neural Network," in Computational Intelligence for Engineering and Management Applications, Singapore, 2023.

[21] B. Roy, L. Malviya, R. Kumar, S. Mal, A. Kumar, T. Bhowmik and J. Hu, "Hybrid Deep Learning Approach for Stress Detection Using Decomposed EEG Signals.," Diagnostics, vol. 13, no. 11, p. 1936, 2023.

[22] A. Ksibi, M. Zakariah, L. Menzli, O. Saidani, L. Almuqren and R. Hanafieh, "Ksibi, A., Zakariah, M., Menzli, L.J., Saidani, O., Almuqren, L. and Hanafieh, R.A.M., Electroencephalography-Based Depression Detection Using Multiple Machine Learning Techniques," Diagnostics, vol. 13, no. 10, p. 1779, 2023.

[23] Q. Abbas, A. Baig and A. Hussain, "Classification of Post-COVID-19 Emotions with Residual-Based Separable Convolution Networks and EEG Signals," Sustainability, vol. 15, no. 2, p. 1293, 2023.

[24] Q. Yao, H. Gu, S. Wang, G. Liang, X. Zhao and X. Li, "Exploring EEG characteristics of multi-level mental stress based on human-machine system. .," Journal of Neural Engineering, vol. 20, no. 5, 2023.

[25] J. Li, D. Lin, Y. Che, J. LV, R. Chen, L. Wang, X. Zeng, J. Ren, H. Zhao and X. Lu, "An innovative EEG-based emotion recognition using a single channel-specific feature from the brain rhythm code method.," Frontiers in Neuroscience, vol. 17, 2023.

[26] P. Anita, D. Chinmayee and A. R. Panat, "Feature Extraction of EEG For Emotion Recognition Using Hjorth Features and Higher Order Crossings," IEEE Xplore, pp. 429-434, 2016.

[27] W. Zheng and B. Lu, "Investigating critical frequency bands and channels for EEG-based emotion recognition with deep neural networks," IEEE Transactions on autonomous mental development, vol. 7, no. 3, pp. 162-175, 2015.

[28] "DEAP: A dataset for emotion analysis using physiological and audiovisual signals," 2022. [Online]. Available: https://www.eecs.qmul.ac.uk/mmv/datasets/deap/. [Accessed 15 March 2022].

[29] A. Baghdadi, Y. Aribi, R. Fourati, N. Halouani, P. Siarry and A. Alimi, "Psychological stimulation for anxious states detection based on EEG-related features," Journal of Ambient Intelligence and Humanized Computing, vol. 12, no. 8, p. 8519–8533, 2021.

[30] M. Islam, A. Rastegarnia and Z. Yang, "Methods for artifact detection and removal from scalp EEG: A review," Neurophysiologie Clinique/Clinical Neurophysiology, vol. 46, no. 4-5, pp. 287-305, 2016.

[31] S. Parihar, P. Shah, R. Sekhar and J. Lagoo, "Model predictive control and its role in biomedical therapeutic automation: A brief review," Applied System Innovation, vol. 5, no. 6, p. 118, 2022.

[32] M. Ladekar, S. Gupta, Y. Joshi and R. Manthalkar, "EEG based visual cognitive workload analysis using multirate IIR filters," Biomedical Signal Processing and Control, vol. 68, p. 102819, 2021.

[33] Q. Gao, A. Omran, Y. Baghersad, O. Mohammadi, M. Alkhafaji, A. Al-Azzawi, S. Al-Khafaji, N. Emami, D. Toghraie and M. Golkar, "Gao, Q., Omran, A.H., Baghersad, Y., Mohammadi, O., Alkhafaji, M.A., Al-Azzawi, A.K.J., Al-Khafaji, S.H., EmamiElectroencephalogram signal classification based on Fourier transform and Pattern Recognition Network for epilepsy diagnosis," Engineering Applications of Artificial Intelligence, vol. 123, p. 106479, 2023.

[34] S. Cohen, O. Katz, D. Presil, O. Arbili and L. Rokach, "Ensemble Learning For Alcoholism Classification Using EEG Signals," IEEE Sensors Journal, vol. 23, no. 5, pp. 1-20, 2023.

[35] X. Wang, Y. Wang, D. Liu, Y. Wang and Z. Wang, "Automated recognition of epilepsy from EEG signals using a combining space–time algorithm of CNN-LSTM," Scientific Reports, vol. 13, no. 1, p. 14876, 2023.

[36] E. Efe and S. Ozsen, "CoSleepNet: Automated sleep staging using a hybrid CNN-LSTM network on imbalanced EEG-EOG datasets," Biomedical Signal Processing and Control, vol. 80, 2023.

[37] W. Cheng, R. Gao, P. Suganthan and K. Yuen, "EEG-based emotion recognition using random Convolutional Neural Networks," Engineering Applications of Artificial Intelligence, vol. 116, p. 105349, 2022.

# Target Detection in Martial Arts Competition Video using Kalman Filter Algorithm Based on Multi target Tracking

Zhiguo Xin

Department of PE, Wuxi Vocational Institute of Arts and Technology, Wuxi, 214200, China

*Abstract*—To solve the low accuracy and poor stability in traditional object tracking methods for martial arts competition videos, a Kalman filtering algorithm based on feature matching and multi object tracking is proposed for object detection in martial arts competition videos. Firstly, feature matching in multi target tracking is studied. Then, based on target feature matching, the Kalman filtering algorithm is fused to construct a target detection model in martial arts videos. Finally, simulation experiments are conducted to verify the performance and application effectiveness of the model. The results showed that the average tracking errors of the model on the X and Y axes were 3.86% and 3.38%, respectively. At the same time, the average accuracy and recall rate in the video target tracking process were 93.64% and 95.48%, respectively. After 100 iterations, the results gradually stabilized. This indicated that the constructed model could accurately detect targets in martial arts competition videos. It had high tracking accuracy and robustness. Compared with traditional object detection methods, this algorithm has better performance and effectiveness. The Kalman filter algorithm based on feature matching and multi target tracking has broad application prospects and research value in target detection in martial arts competition videos.

*Keywords—Multi target tracking; Kalman filtering algorithm; martial arts competition videos; target detection; feature matching*

## I. INTRODUCTION

Object detection is an important research direction in the computer vision, which has broad application value. In martial arts competition videos, object detection can help referees judge the scores of players in real-time, improving the fairness and accuracy of the competition. However, due to the rapid and complex movements in martial arts competitions, traditional object detection methods face a series of challenges in this scenario, such as complex backgrounds and object occlusion [1-3]. In martial arts competitions, multi-target tracking and target detection are key aspects in analyzing the competition process, evaluating athletes' performance, and performing technical and tactical analysis. However, this task is very challenging due to the large number of targets, fast movement speed and frequent occlusion in the scene. Through a large number of studies, it has been found that the utilization of multi-target detection techniques can significantly improve the precision and accuracy of target tracking in martial arts competitions. Therefore, in order to solve the above problems, the study proposes a Kalman filter algorithm based on feature matching and multi-target tracking, which is used for target detection in martial arts competition videos. To address the

above issues, a Kalman filtering algorithm based on feature matching and multi target tracking is proposed for object detection in martial arts competition videos. This algorithm combines two technologies, feature matching and multi target tracking. The target position is accurately located through feature matching. The Kalman filtering algorithm is used to track targets to improve the accuracy and robustness of target detection [4-6]. Firstly, the study focuses on feature matching in multi target tracking. Then, based on target feature matching, the Kalman filtering algorithm is fused to construct a target detection model in martial arts videos. Finally, simulation experiments are conducted to verify the performance and application effectiveness of the model. This model locates the target position through feature matching. The Kalman filtering algorithm is used to track targets to improve the accuracy and robustness of target detection. The research expects to utilize the multi-target tracking Kalman filter algorithm to effectively solve the problems of low target detection accuracy and weak reliability in martial arts competition videos. The contribution of the research is reflected in the utilization of deep learning techniques for feature extraction, which effectively captures the nuances and dynamic changes of targets and improves the accuracy of target detection. Meanwhile, combining the Kalman filter algorithm to predict and correct the target trajectory effectively handles the tracking difficulties caused by target occlusion and fast movement, and enhances the robustness of tracking. By fusing feature extraction, multi-target tracking and Kalman filtering algorithms to construct the model, it can not only focus on the detection and tracking of single targets, but also analyze the interaction and collaborative behaviors of multiple targets, which provides a new perspective for the technical and tactical analysis of martial arts competitions. Compared with the existing techniques, the difference of the research is the organic combination of feature matching and Kalman filtering algorithm. While traditional methods tend to focus only on feature extraction or filter tracking, the research complements the advantages of both to improve the accuracy and robustness of target detection. At the same time, the study also fully considered the characteristics and practical needs of martial arts competitions, making the proposed method more practical and relevant.

Section II is about the related works. In Section III, based on the feature matching algorithm, the Kalman filtering algorithm is fused with it to construct a multi-objective tracking model. Section IV verifies the performance of the

constructed model for comment classification through simulation experiments and practical applications. Finally discussion and conclusion is given in Section V and VI respectively.

## II. RELATED WORKS

Object detection in videos is an important research direction in the computer vision. The main goal is to automatically recognize and track specific target objects in video sequences. This technology has wide applications in many fields, such as racing videos, autonomous driving, security monitoring, medical image analysis, etc. Llano C R et al. State space tracking method based on particle filters for video object tracking. Through experiments, this method had strong performance in tracking objects/people in videos, including foreground/background separation for object movement detection [7]. Lu S and other scholars have investigated the accuracy of real-time video target detection algorithm based on YOLO network for network video image detection. The target information is obtained through image preprocessing and background elimination. Then the convolution operations are applied to reduce the parameters and shorten target detection time. The results showed that this algorithm could significantly shorten the real-time object detection time in videos [8]. Fang Y et al. Multi-intelligent body perception and trajectory prediction method based on spatio-temporal semantics and interaction graph aggregation for effective prediction of spatio-temporal allegories of images and interaction graph aggregation for scene perception and trajectory. The iterative aggregation network was used as background information. Then the trident encoder was decoded and finally detected using prediction methods. The results indicated that this method achieved significant improvements in scene perception and trajectory prediction [9]. Qiu, Ji et al. effectively classified and investigated the performance of small-scale pedestrian detection based on scale prediction method. This method could eliminate the anchor boxes set by most existing detectors. The pixel coordinates of pedestrians at a given center position were predicted. The comprehensive experiments on two real datasets demonstrated that the proposed method achieved excellent performance through [10]. Lyu Y et al. analyzed and studied to improve the detection performance of image detectors in classes without video labels based on agnostic convolutional regression tracker. The performance of the image detector was enhanced through this tracker. This tracker mainly utilized the features of reused image object detectors to learn and track objects. The results indicated that the image detector trained with this tracker could improve accuracy by 5% [11].

Chen Z et al. based on online multi-target tracking algorithm with Kalman filtering and multi-information fusion, conducted a study on leakage detection, false detection and target occlusion during online multi-target tracking. This method utilized Kalman filtering for modeling, and then combined target information with features. The results showed that this method could effectively solve the tracking drift problem caused by target interleaving and occlusion. The main tracking performance parameters were significantly improved [12]. Chen H et al. analyzed and studied the improvement of image target tracking capability based on distributed diffusion traceless Kalman algorithm with covariance intersection strategy. This method could diffuse policy information. Then the adjacent information was fused using a diffusion framework. The results showed that this method could significantly improve the tracking ability of image targets, while also reducing the impact of noise [13]. Liu S et al. based on the improvement algorithm of occlusion prediction tracking based on Kalman filter and spatio-temporal map, the target occlusion, drift and interleaving problems in the target tracking process have been effectively handled and studied. This method could distinguish different images using color histograms and color spatial distribution. The results indicated that the average tracking accuracy of this method was 34.1%. The proposed algorithm improved the performance of multi target tracking process [14].

In summary, the Kalman filtering algorithm is of great significance in the multi target tracking in video images. Based on the Kalman filtering algorithm in image target tracking and detection, the Kalman filtering algorithm can achieve real-time tracking and state estimation of targets in video image target tracking and detection, improving the accuracy and efficiency of target tracking. The research aims to provide a new method for multi object tracking and detection in martial arts competition videos.

## III. DESIGN OF A MARTIAL ARTS COMPETITION VIDEO OBJECT DETECTION MODEL BASED ON FEATURE MATCHING AND KALMAN FILTERING ALGORITHM

Martial arts competition video object detection can provide real-time and accurate object tracking. Through feature matching algorithms, feature extraction and matching can be performed on the characters in the video, thereby accurately tracking the contestants in the competition.

### A. Multiple Target Tracking Algorithm Based on Feature Matching

Athletes' movements in martial arts competitions are fast and complex. Traditional object detection algorithms may not be able to accurately track athletes. The multi target tracking algorithm based on feature matching can effectively analyze and track multiple targets to improve the analysis results of competition videos. In martial arts competitions, players have extremely fast movements. When tracking targets, feature comparison and target matching are performed between the current image and the previous frame image to obtain the correlation of target motion. In the tracking of multi-objective videos, feature extraction and matching of moving targets are required to complete the target tracking. The feature matching of moving targets directly affects the effectiveness of target tracking. Therefore, during feature selection, feature matching is performed on the tracked target to achieve target tracking [15-16]. The selected target matching indicators are the position and height to width ratio of the participants in the rectangular box, as well as the color value of the image. The factors that affect target feature extraction include target area, color, position, and the ratio of height to width. In the detecting the participants, the previous and second frames of images need to be collected. They are first grayed out, and

then subtracted before and after. The difference is binarized before edge detection. During the detection process, pixels are used as the corresponding pixels of the moving target, which represents the position occupied by the target. In the video object tracking, the color of the moving object itself can also serve as a feature matching element. The average color value of the moving target itself is used as a feature for matching. The process of matching and tracking moving targets is shown in Fig. 1.

Based on Fig. 1, to utilize the target feature matching indicators mentioned above, a feature vector is defined to match the target features. This vector can be defined using Eq. (1).

$$a_{n,i} = (s_{n,i}, r_{n,i}, g_{n,i}, b_{n,i}, x_{n,i}, y_{n,i}, rate_{n,i}) \quad (1)$$

In Eq. (1), $a_{n,i}$ represents the feature vector, which is defined as the $i$-th feature vector in the $n$-th frame image. $s_{n,i}$ represents the area occupied by the moving target in the selected image. $r_{n,i}$ represents the average value of red pixels. $g_{n,i}$ represents the average value of green pixels. $b_{n,i}$ represents the average value of blue pixels. $x_{n,i}$ represents the abscissa in the matrix. $y_{n,i}$ represents the ordinate in the matrix. $rate_{n,i}$ represents the ratio of matrix height to width. The variation of the moving target between two frames of images is very small, which makes the image have obvious continuity. It is used as a feature flux to define the similarity function of a target image. Then it is used for feature matching work. The similarity function can be represented by Eq. (2).

$$\Delta s_{i,j} = \frac{\left| s_{n,i} - s_{n-1,j} \right|}{s_{n,j}} \quad (2)$$

In Eq. (2), $\Delta s_{i,j}$ represents the similarity function. $s_{n-1,j}$ represents the area of the $j$-th target in the $n-1$-th frame image. $s_{n,j}$ represents the $j$-th feature vector in the frame image. To determine the color mean of the target, the

similarity function between the three colors in the previous image and the current image is defined. The similarity function of the three colors can be represented by Eq. (3).

$$\begin{cases} \Delta r_{i,j} = \dfrac{\left| r_{n,i} - r_{n-1,j} \right|}{r_{n,i}} \\[2ex] \Delta g_{i,j} = \dfrac{\left| g_{n,i} - g_{n-1,j} \right|}{r_{n,i}} \\[2ex] \Delta b_{i,j} = \dfrac{\left| b_{n,i} - b_{n-1,j} \right|}{b_{n,i}} \end{cases} \quad (3)$$

In Eq. (3), $r_{n-1,j}$ represents the area of the $j$-th target in frame $n-1$ of the red image. $g_{n-1,j}$ represents the area of the $j$-th target in frame $n-1$ of the green image. $b_{n-1,j}$ represents the area of the $j$-th target in frame $n-1$ of the blue image. For the center of the moving target matrix, the similarity function between the current $i$-th target and the $j$-th target in the previous frame image is shown in Eq. (4).

$$\begin{cases} \Delta x_{i,j} = \dfrac{\left| x_{n,i} - x_{n-1,j} \right|}{x_{n,i}} \\[2ex] \Delta y_{i,j} = \dfrac{\left| y_{n,i} - y_{n-1,j} \right|}{y_{n,i}} \end{cases} \quad (4)$$

In Eq. (4), $x_{n-1,j}$ represents the center position of the $j$-th rectangular box in the $n-1$-th frame of the image in the x-axis direction. $y_{n-1,j}$ represents the center position of the $j$-th rectangular box in the $n-1$-th frame of the image in the y-axis direction. For the height to width ratio feature of the bounding rectangle of the moving object in the video, the similarity function is represented by Eq. (5).

$$\Delta rate_{i,j} = \frac{\left| rate_{n,j} - rate_{n-1,j} \right|}{rate_{n,i}} \quad (5)$$



Fig. 1. Flowchart of moving target matching and tracking.

Fig. 2.    Flow chart of multi target tracking algorithm based on feature matching.

In Eq. (5), $rate_{n,j}$ represents the ratio of the height to width of the $j$-th target in the $n$-th frame image. $rate_{n-1,j}$ represents the ratio of the height to width of the $j$-th target rectangle in $n-1$-th frame image. $rate_{n,i}$ represents the ratio of the height to width of the $i$-th target rectangle in the $n$-th frame image. To fuse the four features used for target matching mentioned above together, a metric function is introduced into feature fusion. The definition process of the degree function can be represented by Eq. (6).

In Eq. (6), $a$ is the target feature weighting coefficient. $\beta$ is the color weighting coefficient. $\chi$ represents the x-axis area weighting coefficient. $\gamma$ represents the y-axis target feature weighting coefficient.

$$\Delta a_{i,j} = a(\Delta s_{i,j})^2 + \beta(\Delta r_{i,j})^2 + (\Delta g_{i,j})^2 + \chi(\Delta x_{i,j})^2 + (\Delta y_{i,j})^2 + \gamma(\Delta rate_{i,j})^2 \quad (6)$$

Combined with the analysis of Fig. 2, it can be seen that the feature matching method is used for effective target detection of the features of athletes in the video of the game, and the feature matching process is completed by judging whether it meets the requirements of feature matching by the presence and absence of athletes as well as the size of the threshold value.

*B. Construction of a Multi Target Tracking and Detection ModelBased on Feature Matching and Kalman Filtering Algorithm*

The feature matching algorithm can effectively track the target, but the tracking accuracy is significantly affected if the target is occluded. To address the impact of target occlusion, the Kalman filtering algorithm is introduced on the basis of the feature matching algorithm. The two are fused for the construction of a multi target tracking model. Under the principle of minimum mean square error, a Kalman filter is used to iterate the elements, thereby completing the entire tracking state [17-18]. The fused Kalman filtering algorithm can estimate the past and future states of moving targets based on their current states. The flowchart of the Kalman filtering algorithm after fusion feature matching is shown in Fig. 3.



Fig. 3.    Flowchart of Kalman filtering algorithm.

The tracking of moving target states using Kalman filtering algorithm is affected by random noise. Therefore, the tracking status is first determined. The tracking status is shown in Eq. (7).

$$x_k = Ax_{k-1} + Bu_{k-1} + w_{k-1} \quad (7)$$

In Eq. (7), $A$ represents the state transition matrix. $x_{k-1}$ represents the state $x$ noise value of the $k-1$-th target in the state transition matrix. $B$ represents the state control matrix. $u_{k-1}$ represents the noise value of state $u$ for the $k-1$-th target in the state control matrix. $w_{k-1}$ represents the random noise value. After determining the tracking state, the storage capacity of the tracking state can be found through feature extraction methods. The observation formula is shown in Eq. (8).

$$z_k = Hx_k + v_k \quad (8)$$

In Eq. (8), $H$ represents the observation matrix of the state. $v_k$ represents observation noise. The determination of multi-objective states requires a significant amount of time. Therefore, to simplify the process, the covariance of state noise and observation noise is utilized to reflect the tracking effect by estimating the error in step $k$ of the tracking process. It can be defined by Eq. (9).

$$\begin{cases} \bar{x} = A\bar{x}_{k-1} + Bu_{k-1} \\ p_{\bar{k}} = AP_{k-1} + Q \end{cases} \quad (9)$$

In Eq. (9), $\bar{x}_{k-1}$ represents the tracking result of the previous observation. $\bar{x}$ represents the observation results at the corresponding time. $P_{k-1}$ represents the previous prediction result. $Q$ represents the covariance difference of state noise. After obtaining the tracking status of multiple targets, the observation results are used to determine whether there is an error between the tracking status and the actual observed values. Furthermore, the revised state estimation values and noise values are obtained. It is the process of using the Kalman filtering algorithm to filter the noise. The flowchart of this process is shown in Fig. 4.



Fig. 4. Operation flowchart of Kalman filtering algorithm for noise filtering.

The target motion state in the competition video image is either high-speed or irregular. The common first-order motion model cannot complete the observation of the entire state. Therefore, a second-order motion model is introduced. The Kalman filtering algorithm is used to predict the targets in the second-order motion model to obtain relevant motion features, which are effectively fused with feature extraction algorithms. Assuming that at a certain moment in tracking, the tracked moving target is in a moving rectangle. The velocity of the tracked moving target in the vertical and horizontal directions is uniform motion. Then the motion state needs to meet the uniform motion, as shown in Eq. (10).

$$\begin{cases} v_x(t) = v_x(t-1) \\ v_y(t) = v_y(t-1) \end{cases} \tag{10}$$

In Eq. (10), $v_x(t)$ represents the uniform motion velocity on the x-axis at time $t$. $v_x(t-1)$ represents the uniform motion velocity on the x-axis at time $t-1$. $v_y(t)$ represents the uniform motion velocity on the x-axis at time $t$. $v_y(t-1)$ represents the uniform motion velocity on the y-axis at time $t-1$. The state transition of Kalman filtering can be

determined based on the uniform motion during tracking, as shown in Eq. (11).

$$x_z = Ax_{k-1} + w_{k-1} \tag{11}$$

In Eq. (11), $x_z$ represents the transition state value of the Kalman filter. At this point, the corresponding state transition matrix is shown in Eq. (12).

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \tag{12}$$

The state transition formula can only be used as a directly measured value if it satisfies the matrix. The measurement is shown in Eq. (13).

$$z_c = Hx_k + v_k \tag{13}$$

In Eq. (13), $z_c$ represents the measured value obtained. The corresponding state observation matrix is shown in Eq. (14).

$$H = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix} \tag{14}$$

After analyzing the competition video, there is a significant occlusion of objects in the martial arts competition video. This is mainly caused by the arena, participating athletes, and judges. It is inevitable in the real environment. This occlusion has a significant impact on multi target tracking in competitions. Through the above research, when analyzing the state of moving targets, if occlusion occurs, the image will disappear, and even all images will disappear. If the occluded image merges with the target after a period of time, it can be determined that the image has completely disappeared. If the occluded image appears again in the video rectangle after separating from other targets, it will be used as a new tracking target for tracking and recognition, and matched with a new feature quantity [19-20]. On the basis of image feature matching, the Kalman filtering algorithm is introduced to fuse the two for predicting the motion of moving targets in martial arts competition videos. By utilizing the characteristics of both, the position information of occluded targets is predicted to meet the real-time tracking of targets. The multi target tracking flowchart of this model is shown in Fig. 5.

Fig. 5.    Flow chart of multiple targets tracking in martial arts competition video.

## IV.  Performance Analysis of Multi Target Tracking and Detection Models

To verify the performance of the multi-objective tracking model, 58 videos of different martial arts competitions are obtained through an authorized platform. Each is five minutes. The number of targets in the video varies, ranging from 1 to 3 people. These 58 videos are constructed into a multi target tracking dataset to validate the application performance of the model.

### A.  Performance Analysis of Multi Target Tracking and Detection Models

To analyze the performance of multi target tracking models, the Kalman filtering algorithm and the Minimum Output Sum Square Error (MOSSE) algorithm were compared with the propose method. The error comparison results of the three methods on the X and Y axes of the image are shown in Fig. 6.

In Fig. 6 (a), there was a certain difference in the tracking error of the three methods on the X-axis. The average tracking error of the proposed model was 3.86%. The tracking errors of MOSSE and Kalman methods were 5.94% and 8.05%, respectively. In Fig. 6 (b), the tracking error of the three methods on the Y-axis was smaller than that on the X-axis. The tracking error of the proposed method was 3.38%. The tracking errors of MOSSE and Kalman were 5.17% and 6.23%, respectively. All errors did not exceed 10%. This indicated that the method used to construct the model had high robustness in identifying image targets. To verify the accuracy and recall of the model method in tracking targets, the ratio of the identified targets to the actual targets was used as an evaluation indicator. The comparison results of accuracy and recall were shown in Fig. 7.



(a) Comparison Results of Tracking Errors of Three Methods on the X-axis



(b) Comparison Results of Tracking Errors of Three Methods on the Y-axis

Fig. 6.    Error comparison results of three methods on the X and Y axis of images.

From Fig. 7 (a), all three methods had certain effects in the video target tracking process. The average accuracy of the proposed method in the video target tracking process was 93.64%. The average accuracy of MOSSE and Kalman methods was 81.09% and 78.16%, respectively. Compared to this method, it was 12.55% and 15.48% higher. In Fig. 7 (b), there was also a certain gap in the recall rate of video tracking data among the three methods. The proposed method had a recall rate of 95.48%. The recall rate of MOSSE method was

89.07%. The recall rate of the Kalman method was 83.47%. Recall rate refers to the proportion of correctly predicted positive samples to all actual positive samples. A high recall rate indicates that the model has a stronger ability to correctly predict positive examples. To further validate the multi target tracking performance of the model, three methods were applied to the training and validation sets for comparison. The comparison results were shown in Fig. 8.



(a) Comparison results of average accuracy of three methods          (b) Comparison of recall rates among three methods

Fig. 7.    Comparison results of target tracking accuracy and recall rate in videos using three methods.



(a)Kalman          (b)MOSSE

(c) Model Method

Fig. 8.    Comparison results of loss values between three methods in training and validation sets.

From the comparative analysis in Fig. 8, the proposed method had a faster convergence speed at runtime compared to MOSSE and Kalman methods. After 100 iterations, the results gradually stabilized. At this point, the accuracy difference between the training set and the validation set was very small. The value of the loss function decreased faster and the value was also smaller. This indicated that the proposed method had better stability in the target tracking process compared to the comparison method. Compared to MOSSE and Kalman methods, this method could converge to the exact position of the target faster. The model could learn the features of the target faster, so that it could match more

accurately when the target reappeared. On the training and validation sets, the research model could better fit the features of the target with high accuracy.

### B. The Application Effect of the Multi Target Tracking and Detection Model

To verify the effectiveness of the multi target tracking and detection model in practical applications, the real-time performance of target tracking and the complexity of algorithm operation were used as indicators for verification. The operational efficiency and computational complexity of the three methods were shown in Fig. 9.



(a) Three methods for detecting operational efficiency

(b) Three methods for calculating complexity

Fig. 9.    Comparison of efficiency and computational complexity of three methods for target checking.

In Fig. 9 (a), the operational efficiency of video object detection could reflect the speed of the method in practical applications. The operational efficiency of the proposed method in video object tracking was 77.48%. The operational efficiency of the MOSSE method was 67.34%. The operational efficiency of the Kalman method was 62.55%. In Fig. 9 (b), there were significant differences in the complexity of the three methods in actual operations. When using the proposed method to process images of the same size, the three methods also showed an increasing trend in time consumption as the image size increased. The average time consumption of the proposed method was 0.59ms. The average time consumption of the MOSSE method was 0.85ms. The average time consumption of the Kalman method was 0.93ms. The computational complexity represents the time and space resources required for algorithm execution. Low computational complexity results in less runtime and resources, greatly improving computational efficiency. To further validate the performance of the proposed method, the tracked predicted values were compared with the actual values. The comparison results were shown in Fig. 10.

In Fig. 10, the pixel error range of the true value was between [4.9-13.6]. The pixel error range of the predicted value was between [5.3-12.1]. The difference between the maximum and minimum predicted values and the true values was 1.5 and 0.4. By combining the pixel error trend, the predicted value was basically consistent with the actual trend.

This indicated that the model method had strong applicability in multi target tracking videos. To verify the tracking effect of the model method in multi target video orientation in martial arts competition videos, it was converted into a coordinate system for trajectory prediction. The results of trajectory prediction were shown in Fig. 11.



Fig. 10. Comparison results between video tracking predicted values and true values.

Fig. 11. Comparison of multiple target tracking effects of model methods in coordinate systems.

In Fig. 11, there may be some deviation in the predicted trajectory compared to the actual trajectory. However, the overall trend gap was not significant. Especially in the Y-axis direction, there was a high degree of overlap between the predicted trajectory and the actual trajectory. This may be related to tracking target movement. Despite certain deviations, the predicted trajectory could still roughly reflect the movement trend of the target. This meant that the model could capture the motion patterns of the target and make relatively accurate predictions. To verify the target detection efficiency and multi target tracking ability of the model method, the average detection time and the ability to process multiple targets were used as validation indicators, as shown in Fig. 12.

In Fig. 12 (a), as the number of videos increased, the time required for the proposed method to track the target also showed an upward trend. But the rising speed was not fast, maintaining around 1ms. This indicated that the model method had strong adaptability, which could be well used for target tracking. In Fig. 12 (b), there were a total of 75 moving targets in the entire video image. The proposed method detected a total of 68 moving targets with a detection rate of 91.67%. This indicated that the tracking accuracy and stability of the model method in martial arts competition videos met the design requirements.



(a) Model recognition target time consumption



(b) Multi target tracking performance

Fig. 12. Model method average detection time and multiple target tracking ability results.

## V. DISCUSSION

By analyzing the effective research of multi-target tracking based on feature matching with Kalman filter algorithm for target detection in martial arts competition videos, the method still has some limitations and challenges in performing target detection. First, the data volume and labeling problem is one of the key factors affecting the experimental effect. Due to the relatively small dataset of martial arts competition videos, the training of the deep learning model may not be sufficient, thus affecting the accuracy of target detection. Meanwhile, for the multi-target tracking task, labeling the trajectory of each target is a time-consuming and complex task, especially when there are frequent occlusions and interactions between the targets, and the difficulty of labeling will further increase. Second, the problem of fast target movement and occlusion is also one of the frequently encountered problems in the experiment. In martial arts competitions, the frequent rapid movement of targets with frequent occlusions leads to an increase in the difficulty of feature extraction and also affects the accurate prediction and correction of target trajectories by Kalman filtering. To solve this problem, more advanced target detection algorithms can be tried to improve the accuracy and speed of target detection. In addition, the model generalization ability is also one of the issues that need to be paid attention to in the experiment. Although this study is optimized for martial arts competition videos, the generalization ability of the model may be limited when facing other types of videos or real application scenarios. In order to improve the generalization ability of the model, techniques such as migration learning can be tried to extract more representative features from large-scale datasets to enhance the generalization ability of the model. Finally, real-time performance requirements are also one of the factors to be considered in experiments. For practical applications, such as real-time match analysis or referee assistance systems, the real-time performance of the algorithm is very critical. Therefore, how to improve the running speed of the algorithm while ensuring accuracy is a

problem to be solved. Techniques such as parallel computing can be tried to accelerate the running speed of the algorithm to meet the real-time performance requirements.

In summary, the research on multi-target tracking based on feature matching and Kalman filtering algorithms for target detection in martial arts competition videos faces limitations and challenges in terms of the amount of data, labeling, target motion characteristics, and scene complexity. In order to solve these problems, the performance of the algorithm needs to be further investigated and improved to enhance its target detection accuracy and robustness in martial arts competition videos.

## VI. CONCLUSION

In response to the low accuracy and poor stability in traditional target tracking methods for martial arts competition videos, a multi target tracking and detection model for martial arts competition videos is constructed by integrating feature matching and Kalman filtering algorithms. The results showed that the operational efficiency of the model method in video object tracking was 77.48%, with an average time of 0.59ms. The maximum and minimum predicted values of the proposed method differed from the true values by 1.5 and 0.4, respectively. In the entire video image, there were a total of 75 moving targets. The proposed method detected a total of 68 moving targets with a detection rate of 91.67%. The model performs well in object detection in martial arts competition videos. This model can accurately detect targets in videos under different scenes and lighting conditions. It has high stability and robustness. In addition, the model can also handle the multiple targets to ensure that each target is correctly detected and tracked. Overall, the proposed Kalman filter algorithm based on feature matching and multi target tracking has high accuracy and stability in the target detection model of martial arts competition videos. It can effectively handle the multiple targets, providing an effective technical means for real-time scoring of martial arts competitions. However, there are still shortcomings in the research. There are many complex backgrounds, lighting, etc. in martial arts competitions. In special environments, the tracking and detection capabilities of the proposed method still need to be further improved. Meanwhile, there are deficiencies in the research of multi-target tracking based on feature matching and Kalman filtering algorithm for target detection in martial arts competition videos in terms of data volume and annotation, model generalization ability, and real-time performance requirements. Future research should further expand the dataset, introduce advanced techniques, optimize the algorithm performance, and focus on the requirements of real-time applications to improve the accuracy and robustness of target detection. Interdisciplinary cooperation and communication will also bring new ideas and methods for research in this field.

## REFERENCES

[1] J. Wąsik, D. Mosler, D. Ortenburger, T. Gora, "Stereophotogrammetry measurement of kinematic target effect as speed accuracy benchmark indicator for kicking performance in martial arts," Acta of bioengineering and biomechanics, vol. 23, pp. 1509-1514, 2021.

[2] C. Guo, S. You, Y. Luo, J. Hianlang, W. Dong, "A fast moving target detection tracking and trajectory prediction system for binocular vision," Wuhan University Journal of Natural Sciences, vol. 26, pp. 69-80, 2021.

[3] V. H. Le, "Deep learning-based for human segmentation and tracking, 3D human pose estimation and action recognition on monocular video of MADS dataset," Multimedia Tools and Applications, vol. 82, pp. 20771-20818, 2023.

[4] S. Yang, "Face feature tracking algorithm of aerobics athletes based on Kalman filter and mean shift," International Journal of Biometrics, vol. 14, pp. 394-407, 2022.

[5] R. Golash, Y. K. Jain, "Real-time automatic tracking of hand motion in RGB videos using local feature SIFT." International Journal of Intelligent Systems Design and Computing, vol. 3, pp. 161-177, 2020.

[6] P. Sharma, M. Gangadharappa, "Abnormal behavior detection of stationary objects in surveillance videos with visualization and classification," Concurrency and Computation: Practice and Experience, vol. 34, Pp. 56-70, 2022.

[7] C. R. Llano, Y. Ren, N. I. Shaikh, "Object Detection and Tracking in Real Time Videos," International Journal of Information Systems in the Service Sector (IJISSS), vol. 11, pp. 1-17, 2019.

[8] S. Lu, B. Wang, H. Wang, L. Chen, X. Zhang, "A real-time object detection algorithm for video," Computers & Electrical Engineering, vol. 77, pp. 398-408, 2019.

[9] Y. Fang, B. Luo, T. Zhao, H. Dong, B. B. Jiang, Q. Liu, "ST-SIGMA:Spatio-temporal semantics and interaction graph aggregation for multi-agent perception and trajectory forecasting," CAAI Transactions on Intelligence Technology, vol. 7, pp. 744-757, 2022.

[10] Qiu, Ji, L. Wang, Y. Hu, Y. Wang, "Effective object proposals: size prediction for pedestrian detection in surveillance videos," Electronics Letters 56, vol. 14, pp. 706-709.

[11] Y. Lyu, M. Y. Yang, G. Vosselman, G. S. Xia, "Video object detection with a convolutional regression tracke," ISPRS journal of photogrammetry and remote sensing, vol. 176, pp. 139-150, 2021.

[12] Z. Chen, L. Huang, "Online multi-target tracking algorithm based on Kalman filter and multiple information fusion," Information/Communication, vol. 17, pp. 35-38, 2019.

[13] H. Chen, J. Wang, C. Wang, D. Wang, M. Xin, "Distributed Diffusion Unscented Kalman Filtering Algorithm with Application to Object Tracking," IFAC-PapersOnLine, vol. 53, pp. 3577-3582, 2020.

[14] S. Liu, "The research of multi-object tracking algorithm using Kalman filtering method," International Journal of Innovative Computing and Applications, vol. 10, pp. 107-114, 2019.

[15] Y. Cai, D. Li, X. Wu, R. Song, Y. Wang, "A TLD target tracking algorithm combining CAMSHIFT and Kalman filtering," Computer applications and software, vol. 36, pp. 211-215, 2019.

[16] Y. Lei, "Research on micro video character perception and recognition based on target detection technology," Journal of Computational and Cognitive Engineering, vol. 1, pp. 83-87, 2022.

[17] B. B. Benuwa, B. Ghansah, "Locality-Sensitive Non-Linear Kalman Filter for Target Tracking," International Journal of Distributed Artificial Intelligence (IJDAI), vol. 13, pp. 36-57, 2021.

[18] H. G. Asl, A. D. Firouzabadi, "Integration of the inertial navigation system with consecutive images of a camera by relative position and attitude updating," Proceedings of the Institution of Mechanical Engineers, Part G: Journal of Aerospace Engineering, vol. 233, pp. 5592-5605, 2019.

[19] H. Li, J. Zhu, "Target tracking algorithm based on mean shift and Kalman filter," Journal of Beijing Institute of Technology, vol. 28, pp. 365-370, 2019.

[20] C. Lu, "Kalman tracking algorithm of ping-pong robot based on fuzzy real-time image," Journal of Intelligent & Fuzzy Systems, vol. 38, pp. 3585-3594, 2020.

# 2D-CNN Architecture for Accurate Classification of COVID-19 Related Pneumonia on X-Ray Images

Nurlan Dzhaynakbaev[1], Nurgul Kurmanbekkyzy[2]*, Aigul Baimakhanova[3]*, Iyungul Mussatayeva[4]

Kazakh-Russian Medical University, Almaty, Kazakhstan[1, 2, 3]

Semey Medical University, Semey, Kazakhstan[4]

*Abstract*—In the wake of the COVID-19 pandemic, the use of medical imaging, particularly X-ray radiography, has become integral to the rapid and accurate diagnosis of pneumonia induced by the virus. This research paper introduces a novel two-dimensional Convolutional Neural Network (2D-CNN) architecture specifically tailored for the classification of COVID-19 related pneumonia in X-ray images. Leveraging the advancements in deep learning, our model is designed to distinguish between viral pneumonia, typical of COVID-19, and other types of pneumonia, as well as healthy lung imagery. The architecture of the proposed 2D-CNN is characterized by its depth and a unique layer arrangement, which optimizes feature extraction from X-ray images, thus enhancing the model's diagnostic precision. We trained our model using a substantial dataset comprising thousands of annotated X-ray images, including those of patients diagnosed with COVID-19, patients with other pneumonia types, and individuals with no lung infection. This dataset enabled the model to learn a wide range of radiographic features associated with different lung conditions. Our model demonstrated exceptional performance, achieving high accuracy, sensitivity, and specificity in preliminary tests. The results indicate that our 2D-CNN model not only outperforms existing pneumonia classification models but also provides a valuable tool for healthcare professionals in the early detection and differentiation of COVID-19 related pneumonia. This capability is crucial for prompt and appropriate treatment, potentially reducing the pandemic's burden on healthcare systems. Furthermore, the model's design allows for easy integration into existing medical imaging workflows, offering a practical and efficient solution for frontline medical facilities. Our research contributes to the ongoing efforts to combat COVID-19 by enhancing diagnostic procedures through the application of artificial intelligence in medical imaging.

*Keywords—Machine learning; deep learning; X-Ray; CNN; detection; classification*

## I. INTRODUCTION

The emergence of the COVID-19 pandemic, caused by the SARS-CoV-2 virus, has dramatically reshaped the landscape of healthcare, particularly in the realm of disease diagnosis and management. Manifesting primarily as a respiratory illness often leading to pneumonia, COVID-19 poses unique challenges in terms of rapid and accurate detection, a vital component in controlling the outbreak. In this context, the use of chest X-ray radiography has gained prominence as a key diagnostic tool for COVID-19-related pneumonia, given its accessibility and expediency in comparison to other imaging techniques like CT scans [1]. However, the increased reliance on radiographic analysis has highlighted a need for more automated, efficient, and precise diagnostic methods. Addressing this need, this paper introduces an innovative two-dimensional Convolutional Neural Network (2D-CNN) architecture, designed specifically for classifying X-ray images indicative of COVID-19 related pneumonia [2].

X-ray imaging's pivotal role in detecting COVID-19 related pneumonia is well-documented, offering a rapid and cost-effective means for initial screening [3]. Nonetheless, the interpretation of these images is subject to variability and requires substantial expertise, given the subtlety of early-stage COVID-19 manifestations in the lungs [4]. The advent of deep learning, particularly convolutional neural networks, has shown significant promise in enhancing the accuracy and efficiency of medical image analysis [5]. The 2D-CNN architecture proposed in this research capitalizes on these advancements, focusing on the unique radiographic features of COVID-19 pneumonia, which include peripheral ground-glass opacities and bilateral patchy shadows, distinct from other types of pneumonia and normal lung conditions [6].

Prior research has primarily focused on general pneumonia detection using AI, without special consideration for the specific characteristics of COVID-19 pneumonia [7-8]. Our model is tailored to these unique features, with a deep learning structure that enhances feature extraction and differentiation. An essential aspect of our model is its interpretability, a key factor in medical AI applications, providing clinicians insights into the AI's decision-making process, thereby fostering trust and clinical integration [9]. The training of our model involved a comprehensive dataset of annotated X-ray images, including diverse cases of COVID-19 pneumonia, other pneumonia types, and healthy lungs, ensuring robustness and generalizability [10].

The model's performance in preliminary tests was notable, achieving higher accuracy rates compared to existing pneumonia classification models, a critical factor in clinical settings where diagnostic precision is paramount [11]. False negatives in this context can lead to delayed treatment and increased transmission risk, while false positives can result in unnecessary interventions [12]. Our model's high sensitivity and specificity indicate its potential as a reliable diagnostic aid in the ongoing pandemic [13].

Integration into existing clinical workflows is a crucial factor for the practical application of AI tools in healthcare. The design of our 2D-CNN model facilitates this integration, making it a viable option for rapid deployment in various

healthcare environments, including resource-limited settings [14].

In conclusion, the development of this 2D-CNN architecture for COVID-19 related pneumonia classification marks a significant advancement in medical imaging AI applications. It not only enhances diagnostic accuracy and efficiency but also contributes to global efforts in managing the pandemic. As the healthcare industry continues to navigate the challenges posed by COVID-19, the role of AI becomes increasingly central, and our research underscores the transformative potential of these technologies in medical diagnostics [15].

## II. RELATED WORKS

The application of artificial intelligence (AI) in medical imaging, especially for pneumonia detection, has seen significant advancement in recent years. This section delves into various studies that have contributed to the development of AI in diagnosing respiratory diseases, with a focus on COVID-19-related pneumonia. The early groundwork in applying convolutional neural networks (CNNs) to medical imaging was set by studies [16] and [17], which explored the use of CNNs in detecting common forms of pneumonia from chest X-rays. These foundational studies demonstrated the potential of CNNs to learn complex patterns in medical images, thereby setting the stage for more advanced applications.

With the onset of the COVID-19 pandemic, the focus shifted towards differentiating COVID-19 pneumonia from other types. Research in [18] and [19] specifically targeted the development of deep learning models trained on COVID-19 X-ray datasets. These studies were pivotal in identifying the unique radiographic features of COVID-19, such as ground-glass opacities and bilateral infiltrates, and how AI models could be trained to recognize these features with high accuracy. The challenge of dataset diversity and size was addressed in studies [20], [21], and [22], emphasizing the importance of comprehensive datasets in developing robust CNN models. These works discussed how a diverse range of X-ray images, including data augmentation techniques, could enhance the model's ability to generalize across different presentations of COVID-19.

Transfer learning emerged as a significant technique in medical imaging AI, as highlighted in research [23] and [24]. These studies successfully adapted pre-trained models from non-medical domains to medical datasets, demonstrating the effectiveness of this approach in rapidly deploying AI solutions for emerging health crises like COVID-19. Further extending the capabilities of CNNs, studies [25] and [26] focused on grading the severity of lung infections. This approach went beyond mere detection and provided critical insights into the extent of lung involvement, which is crucial for treatment planning in COVID-19 cases.

The interpretability of AI models in medical diagnosis gained attention in studies [27] and [28]. These works introduced methods for visualizing the decision-making process of CNNs, which is vital for gaining the trust of clinicians in AI-assisted diagnoses. Comparative studies, such as those in [29] and [30], evaluated various CNN architectures

to determine the most effective models for medical image analysis. The insights from these comparisons have been instrumental in guiding the development of efficient and accurate models for pneumonia detection.

The integration of AI models into clinical workflows was explored in studies [31] and [32]. These works examined the practical aspects of implementing AI tools in healthcare settings, emphasizing the need for user-friendly, practical models for medical professionals. Specific AI architectures were the focus of research [33] and [34], which delved into optimizing layer structures in CNNs for better feature extraction from medical images. These findings have informed the development of sophisticated AI models capable of detecting subtle anomalies in X-ray images.

Ensemble methods, combining multiple AI models for improved diagnostic accuracy, were explored in studies [35] and [36]. These approaches showed potential in reducing misdiagnosis by leveraging the strengths of different AI architectures. An integrated approach using clinical data alongside imaging data in AI models was presented in research [37] and [38]. This holistic method resulted in more nuanced and accurate diagnoses by considering both radiographic features and patient history. Studies in [39] and [40] addressed the scalability and adaptability of AI models, particularly in resource-limited settings. These works emphasized the need for accessible and effective AI solutions in diverse healthcare environments. The interdisciplinary application of natural language processing (NLP) in medical imaging was explored in studies [41] and [42]. These approaches utilized NLP to extract information from radiology reports, providing additional context to AI models and enhancing diagnostic accuracy. Ethical considerations in the deployment of AI in healthcare were discussed in studies [43] and [44]. These studies focused on responsible AI use, patient privacy, and addressing potential biases in AI models.

Finally, future directions in medical imaging AI, as speculated in [45] and [46], include integrating AI models with other diagnostic tools and evolving AI algorithms to adapt to the changing landscape of diseases like COVID-19.

In conclusion, the body of work from [16] to [47] provides a comprehensive overview of the advancements and challenges in applying AI to the diagnosis of pneumonia and respiratory diseases. These studies underscore the potential of AI to revolutionize medical imaging, offering enhanced patient care and management, especially in response to global health crises like the COVID-19 pandemic.

## III. MATERIALS AND METHODS

The Materials and Methods section is a cornerstone of any scientific research paper, providing the necessary details to understand and replicate the study. In this section, we outline the comprehensive approach undertaken in our research, which involves the development and validation of a two-dimensional Convolutional Neural Network (2D-CNN) architecture for the classification of COVID-19 related pneumonia in X-ray images. This section is structured to detail the dataset selection and preparation, the design and implementation of the 2D-CNN model, and the methodologies employed for training,

testing, and validating the model. Additionally, we describe the statistical methods used for the analysis of the results, ensuring a transparent and reproducible research process. Fig. 1 demonstrates an example of a pneumonia caused by COVID-19.

### A. Data

The "Covid19-Pneumonia-Normal Chest X-Ray Images" dataset serves as an invaluable resource for researchers and the medical community, particularly in the domain of applying deep learning techniques for the detection and classification of COVID-19 and pneumonia from chest X-ray images.

This dataset is efficiently organized into three distinct subfolders, namely COVID, NORMAL, and PNEUMONIA, each containing chest X-ray (CXR) images corresponding to their respective classifications. Such a structure facilitates easy access and processing of the data for research purposes. Specifically, the dataset comprises 1,626 images of COVID-19 cases, 1,802 images of normal cases, and 1,800 images of pneumonia cases. This comprehensive collection allows for a balanced representation of each category, which is crucial for training and validating deep learning models with a high degree of accuracy.

A notable feature of this dataset is the standardization of all images. Each image has been preprocessed and resized to a uniform dimension of 256x256 pixels in PNG format. This uniformity is vital for maintaining consistency across the dataset, ensuring that the deep learning models can learn and classify the images without bias towards different sizes or formats.

The significance of this dataset extends beyond its structural organization and preprocessing. It provides a critical tool for advancing research in medical imaging, especially in the current global health context where rapid and accurate diagnosis of COVID-19 is essential. By offering a substantial number of categorized images, it enables the development of sophisticated AI models capable of distinguishing between COVID-19, pneumonia, and normal chest conditions with high precision.

Researchers utilizing this dataset are encouraged to cite pertinent articles that have contributed to its development. Key references include publications by [47-48] These works delve into the architecture and effectiveness of deep convolutional neural networks for classifying COVID-19 in chest X-ray images, providing a foundation for further research in this field.

In conclusion, the "Covid19-Pneumonia-Normal Chest X-Ray Images" dataset is a vital asset for the ongoing development of AI in medical diagnostics, particularly for classifying various lung conditions in the era of COVID-19. Its comprehensive, well-organized, and standardized collection of images is instrumental for researchers striving to enhance the accuracy and efficiency of diagnostic methods through deep learning techniques. Fig. 2 demonstrates samples of the applied dataset.

Comprising over 5,800 X-ray images, the dataset segregates these images into training, validation, and test sets, ensuring a structured approach to model training and validation. Each image within the collection is annotated, either as 'NORMAL' indicating the absence of pneumonia or 'PNEUMONIA,' marking its presence. Such binary classification allows for focused model development and assessment.

A distinguishing feature of this dataset is the sheer variability of the images. Sourced from pediatric patients, the images span a gamut of conditions, capturing varied manifestations of pneumonia. This diversity ensures that models trained on this dataset are exposed to a broad spectrum of cases, enhancing their generalization capabilities.



Fig. 1. Chest pneumonia explanation.

Fig. 2. Samples of the applied dataset.

Fig. 3 demonstrates distribution of classes in the applied dataset. From the figure, we can observe that the number of images in each category is relatively balanced, with a slight variation in the counts. Such a distribution is beneficial for training machine learning models, as it provides a diverse set of examples for each class, helping the model to learn and generalize better across different conditions.

The displayed images provide a visual overview of the three categories—COVID, NORMAL, and PNEUMONIA—in the applied dataset. For each category, a few sample images have been randomly selected to illustrate the typical characteristics visible in chest X-rays.

Fig. 3.   Samples of normal and pneumonia chest X-Rays.

COVID: The images in this category are from patients diagnosed with COVID-19. Radiographic features specific to COVID-19, such as ground-glass opacities and bilateral infiltrates, might be observable in these X-rays.

NORMAL: These images represent normal chest X-rays from individuals without lung infections. They typically exhibit clear lung fields without the opacities or infiltrates seen in the other two categories.

PNEUMONIA: The images here are from patients with various forms of pneumonia, other than COVID-19. Pneumonia X-rays often show areas of increased opacity, which can be localized or diffuse, depending on the type and severity of the infection.

*B. Proposed Model*

In the critical field of medical diagnostics, where precision is of utmost importance, the described sequential model presents a sophisticated computational framework specifically engineered for the detection of pneumonia through medical imaging. This innovative model adeptly merges the capabilities of a well-established pre-trained architecture with tailor-made layers, thereby enabling the meticulous extraction of intricate features and facilitating effective classification. The architecture of this pioneering deep learning model for pneumonia classification is depicted in Fig. 4.

At the core of this model lies the VGG16 layer (Functional), a convolutional neural network originally developed by the Visual Geometry Group at the University of Oxford. This pre-trained layer, comprising 14,714,688 parameters, excels in extracting complex, hierarchical features from the input imagery. It outputs a tensor of dimensions $8\times8\times512$, which represents a rich set of features extracted from the X-ray images. These features are essential for the nuanced detection of pneumonia, underscoring the VGG16 layer's critical role in the model's architecture.

In the advanced architecture following the VGG16 layer, the model incorporates a flatten layer. This layer plays a pivotal role in transforming the three-dimensional feature tensor output from the previous layer into a one-dimensional vector. This transformation is a critical step, enabling the seamless integration of the convolutional output with

subsequent dense layers, effectively linking the feature extraction process to the classification phase.

To address the prevalent issue of overfitting, where models perform well on training data but underperform with new, unseen data, the model includes a dropout layer. This layer enhances the model's generalization capabilities by randomly deactivating a subset of neurons during training epochs. This introduction of randomness fosters a level of robustness within the model, ensuring more consistent performance across various datasets.

Following this, a dense layer comprising 4,194,432 parameters and 128 neurons is integrated. This fully connected layer acts as an intermediary, processing the one-dimensional flattened features derived from the preceding layers. This stage is instrumental in the progressive journey of classification within the model.

To further reinforce the model's robustness and mitigate overfitting, a secondary dropout layer is employed. This layer re-emphasizes the model's commitment to regularization, bolstering its ability to generalize across different datasets.



Fig. 4.   Proposed 2D-CNN architecture.

The architecture culminates with a final dense layer consisting of two neurons. With only 258 parameters, this layer is responsible for the concluding step of the classification process. It outputs probabilistic scores for the two classes—'NORMAL' and 'PNEUMONIA', thereby finalizing the model's decision-making pathway in distinguishing between the two categories.

*C. Evaluation Parameters*

In this research, the concept of accuracy emerges as a critical metric for evaluating the performance of the developed deep learning model in classifying chest X-ray images into COVID-19, pneumonia, and normal categories. Accuracy, in this context, is defined as the proportion of correctly classified images out of the total number of images evaluated. This measurement encapsulates the model's effectiveness in correctly identifying each class and is a fundamental indicator of its diagnostic reliability. High accuracy is essential in medical diagnostics, as it directly impacts the quality of patient care and treatment decisions. An accurate model ensures confidence in automated diagnoses, reducing the likelihood of misdiagnosis and subsequently enhancing patient outcomes. Throughout the research, maintaining and improving the accuracy of the model has been a primary focus, with the goal of developing a tool that is not only technologically advanced but also clinically dependable and effective in real-world healthcare settings [49].

$$accuracy = \frac{TP + TN}{TP + FN + TN + FP}, \tag{1}$$

In the realm of this research, precision is an indispensable metric, pivotal in assessing the model's capability to classify chest X-ray images into distinct categories of COVID-19, pneumonia, and normal. Precision, specifically, refers to the proportion of true positive predictions relative to the total number of positive predictions made by the model. It is a measure of the model's ability to correctly identify positive instances among all instances it labeled as positive. In a clinical setting, high precision is crucial as it minimizes the rate of false positives – instances where the model incorrectly identifies a condition. Especially in medical diagnostics, this is vital, as false positives can lead to unnecessary anxiety, further testing, and potentially unwarranted treatment. Therefore, in developing the deep learning model, a significant emphasis is placed on enhancing precision, ensuring that the model not only identifies conditions accurately but also limits the occurrence of false alarms, thus providing reliable and trustworthy diagnostic support [50].

$$precision = \frac{TP}{TP + FP}, \tag{2}$$

In this research, recall, also known as sensitivity, is a key performance metric for the deep learning model developed for classifying chest X-ray images into categories like COVID-19, pneumonia, and normal. Recall is defined as the proportion of actual positive cases correctly identified by the model, essentially measuring the model's ability to detect true positives from all the actual positive cases. In the context of

medical diagnostics, a high recall rate is extremely important, as it reflects the model's effectiveness in identifying all relevant instances of a condition, thereby reducing the risk of missing a diagnosis. This is particularly crucial for conditions like COVID-19 and pneumonia, where timely and accurate detection can significantly impact patient outcomes and treatment decisions. Therefore, optimizing recall in the model ensures that it not only identifies conditions accurately but also minimizes false negatives, making it a reliable tool in detecting cases that require immediate medical attention.

$$recall = \frac{TP}{TP + FN}, \tag{3}$$

In this research, the F-score (or F1 score) serves as a crucial statistical measure to gauge the precision and recall balance of the developed deep learning model for classifying chest X-ray images into COVID-19, pneumonia, and normal categories. The F-score is the harmonic mean of precision and recall, providing a single metric that encapsulates both the accuracy of the model's positive predictions and its ability to identify all relevant instances. This metric is particularly valuable in medical diagnostics, where it is essential to strike a balance between not missing actual cases (high recall) and minimizing false alarms (high precision). The relevance of the F-score in this context lies in its ability to provide a comprehensive view of the model's performance, especially in scenarios where an equal trade-off between precision and recall is desired. In summary, the F-score is an integral part of evaluating the model's efficacy, ensuring that it is not only accurate but also reliable in practical clinical applications.

$$F - score = \frac{2 \times precision \times recall}{precision + recall}, \tag{4}$$

In this research, the Receiver Operating Characteristic (ROC) curve and the Area Under the Curve (AUC) are pivotal in evaluating the performance of the deep learning model designed for classifying chest X-ray images into COVID-19, pneumonia, and normal categories. The ROC curve is a graphical representation that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. It plots the True Positive Rate (Recall) against the False Positive Rate, providing insight into the trade-off between sensitivity and specificity at various threshold levels. The AUC, a key component of this analysis, quantifies the entire two-dimensional area underneath the entire ROC curve. A higher AUC value indicates better model performance, with a value of 1 representing a perfect classifier. In the context of medical imaging, ROC-AUC is particularly crucial as it encompasses the model's overall capability to distinguish between the classes across all possible thresholds, offering a robust measure of its diagnostic accuracy.

## IV. EXPERIMENTAL RESULTS

In the Experiment Results section of this research, we delve into the empirical findings derived from the application of our deep learning model to the task of classifying chest X-ray images into COVID-19, pneumonia, and normal categories. This section meticulously presents the outcomes of various

tests conducted to evaluate the model's performance, offering a detailed analysis of its effectiveness and reliability. Key performance metrics such as accuracy, precision, recall, F-score, and ROC-AUC are thoroughly examined, providing a comprehensive understanding of the model's capabilities [51]. The results are contextualized with comparative analyses and discussions, shedding light on the model's strengths and areas for improvement. This section not only validates the model's efficacy through quantitative measures but also offers critical insights into its practical applicability in the realm of medical diagnostics. By exploring the experimental results, we aim to underscore the significance of our model in enhancing diagnostic accuracy and contributing to the advancement of AI applications in healthcare.

Fig. 5 demonstrates classification results of x-rays using the proposed 1D CNN model. Model classifies the input images into three types as normal X-rays, COVID, and Pneumonia cases.

Fig. 6 presented showcases a remarkable instance of the deep learning model's capability in detecting lung opacity in a chest X-ray image. It vividly illustrates the area where lung opacity is identified, highlighted by the model's advanced image processing and feature detection algorithms. This visual representation is a clear demonstration of the model's precision in pinpointing areas of concern within the lung, a crucial aspect in diagnosing conditions such as pneumonia or COVID-19. The highlighted region in the X-ray image specifically denotes the detected opacity, offering a clear and concise visual cue for medical professionals. This figure not only exemplifies the model's diagnostic accuracy but also underscores its potential as a valuable tool in aiding clinicians in the rapid and effective assessment of pulmonary conditions. The clarity and precision of the image highlight the advancements in AI-driven medical imaging and its growing significance in enhancing diagnostic processes.



Fig. 5. Obtained classification results.

Fig. 6. Lung opacity detection results.

In the labyrinth of scientific inquiry, the Results section emerges as a lighthouse, illuminating the empirical achievements and performance indicators attained in our investigative endeavor. Grounded in a foundation of thorough experimentation and detailed data examination, the forthcoming results elucidate the effectiveness and implications of our employed methodologies. This segment endeavors to provide a clear and exhaustive depiction of the model's proficiency in classifying pneumonia through X-ray imagery, evaluated against established performance metrics. As we navigate through the detailed terrains of accuracy, precision, recall, among other assessment criteria, we invite our readers to assess the strengths and limitations inherent in our investigative method. We now embark on this analytical voyage, transitioning from raw data to enlightening discoveries.

Fig. 7 presents a graphical depiction of the training and validation accuracy achieved over 25 learning epochs. The model demonstrates commendable learning efficiency, achieving a notable accuracy of 90% in the initial epochs. This level of accuracy is further enhanced as the learning progresses. By the end of the 25th epoch, the model reaches a peak accuracy of 96%, underscoring its potent learning capacity and the robustness of its architecture in effectively classifying X-ray images.

Fig. 8 offers a graphical elucidation of the training and validation losses encountered throughout 25 learning epochs. The trajectory of the training loss is delineated by a blue line, while the validation loss is represented by a red line. A close examination of this figure shows a consistent diminution in both training and validation losses from the commencement of the learning cycle. This trend is indicative of the model's effective learning and adaptation capabilities as it progresses through each epoch. Upon reaching the termination of the 25 epochs, a convergence of both loss metrics is observed, with each arriving at their respective lowest points. This convergence is a testament to the model's proficiency in minimizing the divergence between predicted outcomes and actual data. The observed pattern in the losses is reflective of a model that has undergone substantial training and refinement, reaching a state of maturity by the end of the designated learning epochs.

Fig. 7.   Model accuracy results.



Fig. 8.   Model loss results.

Fig. 9 demonstrates the Receiver Operating Characteristic (ROC) curve for the proposed model, with a value of 0.8, provides a significant insight into its diagnostic ability, particularly in distinguishing between different classes in the classification task. An ROC curve is a graphical representation that plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings. The curve essentially evaluates the trade-offs between sensitivity (true positive rate) and specificity (1 - false positive rate).

In the context of this model, an ROC value of 0.8 indicates a high level of discriminative ability. This means the model has a strong capacity to correctly identify true positives while minimizing false positives. An ROC value of 1.0 would represent a perfect model with 100% sensitivity and specificity,

while a value of 0.5 would suggest no discriminative ability better than random chance.

A value of 0.8 suggests that the model effectively balances sensitivity and specificity. This is particularly important in medical diagnostics, where the ability to correctly identify true cases (sensitivity) without wrongly labeling negative cases as positive (specificity) is crucial. In practical terms, this level of ROC indicates that the model is quite reliable in its classifications, though there is still room for improvement to reach near-perfect classification accuracy.

In summary, the ROC curve with a value of 0.8 for the proposed model is indicative of its robust performance in classifying chest X-ray images, striking a commendable balance between identifying true cases of the condition and avoiding false alarms.

Fig. 9.    ROC-AUC curve results.

As we conclude, it is clear that the presented findings offer a comprehensive understanding of the proposed model's performance in classifying chest X-ray images. The metrics discussed, including accuracy, precision, recall, F-score, ROC-AUC, and the analysis of training and validation losses, collectively paint a picture of a robust and effective model. While the results are promising, indicating a high degree of reliability and efficiency, they also pave the way for further enhancements and explorations.

The journey through these experimental results has been illuminating, revealing both the strengths and potential areas for improvement in our model. It is evident that the field of medical image classification, particularly in the realm of pneumonia detection, is fertile ground for continued research and development.

In moving forward, these results will serve as a foundation for future work, guiding refinements in the model and inspiring new approaches to enhance its accuracy and usability in clinical settings. The insights gained here are not only valuable for the specific task of pneumonia classification but also contribute to the broader discourse in the application of AI in healthcare diagnostics. With these conclusions, we close this section, carrying forward the knowledge and understanding gleaned to inform subsequent phases of our research endeavor.

## V.    DISCUSSION

The Discussion section of this paper provides a comprehensive analysis of the findings from the experiment results, offering a deeper insight into the implications, limitations, and potential future directions of the research. The proposed model, leveraging a deep learning approach for the classification of pneumonia in chest X-ray images, demonstrates promising results, which aligns with the findings in recent studies [52]. However, a critical evaluation of these results, in light of existing literature and emerging trends in medical imaging, is imperative for a holistic understanding.

### A.    Model Performance and Comparison with Existing Methods

The high accuracy and ROC-AUC score achieved by our model are significant accomplishments, underscoring its potential as a reliable tool in medical diagnostics. This aligns with the trends observed in similar studies [53], where deep learning models have shown considerable success in medical image analysis. The precision and recall metrics also indicate a balanced model performance, essential in medical applications to reduce both false positives and false negatives. However, when compared to similar models in the literature [54], it is evident that while our model performs commendably, there is still room for improvement, especially in terms of achieving consistency across various datasets.

### B.    Importance of Dataset Quality and Diversity

The dataset's quality and diversity played a pivotal role in the model's performance, a finding consistent with observations made in studies [55]. The diverse range of images in the dataset helped in training a more robust model, capable of generalizing across different patient demographics and image qualities. This reinforces the notion that for deep learning models, especially in medical imaging, the dataset's comprehensiveness and representativeness are as crucial as the model architecture itself.

### C.    Impact of Overfitting and Regularization Techniques

The incorporation of dropout layers to combat overfitting proved to be effective, as reflected in the convergence of training and validation losses. This approach aligns with the strategies recommended in recent research [56], emphasizing the importance of regularization techniques in improving model generalizability. However, it's worth noting that while dropout layers aid in reducing overfitting, they are not a panacea, and continuous monitoring of model performance is necessary to ensure its reliability.

### D.    Implications in Clinical Settings

The application of this model in clinical settings holds significant promise. Its ability to accurately classify pneumonia from chest X-rays can aid in quicker diagnosis and treatment, potentially improving patient outcomes. However, as suggested in previous studies [49], the integration of AI tools in clinical practice requires careful consideration of factors like user acceptance, interpretability, and alignment with clinical workflows.

### E.    Limitations and Future Directions

One of the primary limitations of this study is the dependency on the quality and diversity of the dataset. As shown in previous research [52], models trained on limited or biased datasets can exhibit reduced performance in real-world scenarios. Future work should focus on expanding the dataset to include a wider variety of images, including those from underrepresented groups and varied clinical settings.

Another area for future exploration is the interpretability of the model. As AI applications in healthcare continue to grow,

the need for models that provide not just accurate, but also interpretable results becomes crucial [53]. Developing techniques that offer insights into the model's decision-making process could enhance clinician trust and aid in the broader acceptance of AI tools in medical diagnostics.

*F. Contribution to the Field*

This research contributes to the growing body of work on the application of deep learning in medical imaging. By offering a model that demonstrates high accuracy and robustness in pneumonia classification, it paves the way for further advancements in this field. Additionally, the insights gained from this study regarding dataset importance, model generalizability, and the challenges of integrating AI into clinical practice provide valuable guidance for future research endeavors.

In conclusion, the findings from this research offer promising prospects for the use of deep learning in medical image analysis, particularly in the classification of pneumonia from chest X-rays. While the results are encouraging, continuous efforts in refining the model, expanding the dataset, and enhancing interpretability are essential for the successful translation of these findings into clinical practice. This research not only contributes to the technological advancements in medical diagnostics but also highlights the critical considerations necessary for the effective and ethical application of AI in healthcare.

## VI. CONCLUSION

In concluding this research paper, it is imperative to reflect on the significant strides made in the realm of medical diagnostics through the application of advanced deep learning techniques. The development and implementation of a convolutional neural network (CNN) model for the classification of pneumonia from chest X-ray images represent a notable advancement in the field. The model's ability to accurately distinguish between COVID-19, pneumonia, and normal cases, as evidenced by the high accuracy, precision, recall, and ROC-AUC scores, underscores the potential of AI in enhancing diagnostic processes. The robust performance of the model, facilitated by the comprehensive and diverse dataset, as well as effective regularization techniques to counter overfitting, marks a crucial step towards the integration of AI in clinical settings. However, it is essential to acknowledge the limitations encountered, particularly the dependency on dataset quality and the challenges of ensuring model interpretability and integration within clinical workflows.

Looking ahead, this research paves the way for future explorations in medical image analysis using AI. The insights gained underscore the need for ongoing efforts to expand and diversify training datasets to enhance the model's applicability and reliability across varied clinical scenarios. Additionally, the quest for improved interpretability of AI models remains paramount, as this is crucial for clinician acceptance and ethical deployment in healthcare settings. The integration of AI tools like the proposed model in clinical practice requires a multifaceted approach, involving not only technological advancements but also considerations of user experience, workflow compatibility, and ethical implications. In summary, this research contributes significantly to the field of AI in medical diagnostics, offering a promising tool for pneumonia classification while also highlighting the critical areas for continued research and development. The journey of integrating AI in healthcare is ongoing, and the findings from this study provide valuable guidance and impetus for future advancements in this dynamic and impactful field.

## REFERENCES

[1] Tursynova, A., & Omarov, B. (2021, November). 3D U-Net for brain stroke lesion segmentation on ISLES 2018 dataset. In 2021 16th International Conference on Electronics Computer and Computation (ICECCO) (pp. 1-4). IEEE.

[2] Farhan, A. M. Q., & Yang, S. (2023). Automatic lung disease classification from the chest X-ray images using hybrid deep learning algorithm. Multimedia Tools and Applications, 1-27.

[3] Aslani, S., & Jacob, J. (2023). Utilisation of deep learning for COVID-19 diagnosis. Clinical Radiology, 78(2), 150-157.

[4] Patel, B., & Thakkar, M. M. (2023). Deep Learning Approaches in Pediatric Pulmonary Diagnostics: Evaluating 1D-CNN, 2D-CNN, and 3D-CNN Variants on Chest X-Ray Classifications. Journal of Korean Academy of Psychiatric and Mental Health Nursing, 5(4), 163-173.

[5] Ahmed, M., Gilanie, G., Ahsan, M., Ullah, H., & Sheikh, F. A. (2023). Review of Artificial Intelligence-based COVID-19 Detection and A CNN-based Model to Detect Covid-19 from X-Rays and CT images. VFAST Transactions on Software Engineering, 11(2), 100-112.

[6] UmaMaheswaran, S. K., Prasad, G., Omarov, B., Abdul-Zahra, D. S., Vashistha, P., Pant, B., & Kaliyaperumal, K. (2022). Major challenges and future approaches in the employment of blockchain and machine learning techniques in the health and medicine. Security and Communication Networks, 2022.

[7] Preethi, M. S., Rajitha, B., Reddy, K. S., Kovela, B., & Veeramalla, S. K. (2024). COVID-19 Detection Using Convolutional Neural Networks from Chest X-Ray Images. In Internet of Medical Things in Smart Healthcare (pp. 11-37). Apple Academic Press.

[8] Sarkar, O., Islam, M. R., Syfullah, M. K., Islam, M. T., Ahamed, M. F., Ahsan, M., & Haider, J. (2023). Multi-Scale CNN: An Explainable AI-Integrated Unique Deep Learning Framework for Lung-Affected Disease Classification. Technologies, 11(5), 134.

[9] Aboamer, M. A., Sikkandar, M. Y., Gupta, S., Vives, L., Joshi, K., Omarov, B., & Singh, S. K. (2022). An investigation in analyzing the food quality well-being for lung cancer using blockchain through cnn. Journal of Food Quality, 2022.

[10] Bayoudh, K., Hamdaoui, F., & Mtibaa, A. (2020). Hybrid-COVID: a novel hybrid 2D/3D CNN based on cross-domain adaptation approach for COVID-19 screening from chest X-ray images. Physical and engineering sciences in medicine, 43, 1415-1431.

[11] Lee, M. H., Shomanov, A., Kudaibergenova, M., & Viderman, D. (2023). Deep Learning Methods for Interpretation of Pulmonary CT and X-ray Images in Patients with COVID-19-Related Lung Involvement: A Systematic Review. Journal of Clinical Medicine, 12(10), 3446.

[12] Qi, Q., Qi, S., Wu, Y., Li, C., Tian, B., Xia, S., ... & Yu, H. (2022). Fully automatic pipeline of convolutional neural networks and capsule networks to distinguish COVID-19 from community-acquired pneumonia via CT images. Computers in Biology and Medicine, 141, 105182.

[13] Thai-Nghe, N., Hong, N. M., Nhu, P. T. B., & Hai, N. T. (2023, October). Classification of Pneumonia on Chest X-ray Images Using Transfer Learning. In International Conference on Intelligence of Things (pp. 85-93). Cham: Springer Nature Switzerland.

[14] Goncharov, M., Pisov, M., Shevtsov, A., Shirokikh, B., Kurmukov, A., Blokhin, I., ... & Belyaev, M. (2021). CT-Based COVID-19 triage: Deep multitask learning improves joint identification and severity quantification. Medical image analysis, 71, 102054.

[15] Podder, P., Das, S. R., Mondal, M. R. H., Bharati, S., Maliha, A., Hasan, M. J., & Piltan, F. (2023). Lddnet: a deep learning framework for the diagnosis of infectious lung diseases. Sensors, 23(1), 480.

[16] Majhi, V., & Paul, S. (2023). Comparative Analysis on Available Technique for the Detection of Covid-19 through CT-Scan and X-Ray using Machine Learning: A Systematic Review.

[17] Serena Low, W. C., Chuah, J. H., Tee, C. A. T., Anis, S., Shoaib, M. A., Faisal, A., ... & Lai, K. W. (2021). An overview of deep learning techniques on chest X-ray and CT scan identification of COVID-19. Computational and Mathematical Methods in Medicine, 2021, 1-17.

[18] Saju, B., Tressa, N., Dhanaraj, R. K., Tharewal, S., Mathew, J. C., & Pelusi, D. (2023). Effective multi-class lungdisease classification using the hybridfeature engineering mechanism. Mathematical Biosciences and Engineering, 20(11), 20245-20273.

[19] Modak, S., Abdel-Raheem, E., & Rueda, L. (2023). Applications of deep learning in disease diagnosis of chest radiographs: A survey on materials and methods. Biomedical Engineering Advances, 100076.

[20] Liu, W., Ni, Z., Chen, Q., & Ni, L. (2023). Attention-guided Partial Domain Adaptation for Automated Pneumonia Diagnosis From Chest X-Ray Images. IEEE Journal of Biomedical and Health Informatics.

[21] Das, A. K., Kalam, S., Kumar, C., & Sinha, D. (2021). TLCoV-An automated Covid-19 screening model using Transfer Learning from chest X-ray images. Chaos, Solitons & Fractals, 144, 110713.

[22] Wu, Y., Qi, Q., Qi, S., Yang, L., Wang, H., Yu, H., ... & Chen, R. (2023). Classification of COVID-19 from community-acquired pneumonia: Boosting the performance with capsule network and maximum intensity projection image of CT scans. Computers in Biology and Medicine, 154, 106567.

[23] Narin, A., Kaya, C., & Pamuk, Z. (2021). Automatic detection of coronavirus disease (covid-19) using x-ray images and deep convolutional neural networks. Pattern Analysis and Applications, 24, 1207-1220.

[24] Kirana, K. C., Wibawanto, S., & Hamdan, A. (2022). Optimization of 2D-CNN Setting for the classification of covid disease using Lung CT Scan. Jurnal Ilmu Komputer dan Informasi, 15(2), 143-149.

[25] Panigrahi, S., Raju, U. S. N., Pathak, D., Kadambari, K. V., & Ala, H. (2023). Rapid detection of COVID-19 from chest X-ray images using deep convolutional neural networks. International Journal of Biomedical Engineering and Technology, 41(1), 1-15.

[26] Kaur, J., & Kaur, P. (2022). Outbreak COVID-19 in medical image processing using deep learning: a state-of-the-art review. Archives of Computational Methods in Engineering, 1-32.

[27] Santosh, K. C., GhoshRoy, D., & Nakarmi, S. (2023, August). A systematic review on deep structured learning for COVID-19 screening using chest CT from 2020 to 2022. In Healthcare (Vol. 11, No. 17, p. 2388). MDPI.

[28] Qi, S., Xu, C., Li, C., Tian, B., Xia, S., Ren, J., ... & Yu, H. (2021). DR-MIL: deep represented multiple instance learning distinguishes COVID-19 from community-acquired pneumonia in CT images. Computer Methods and Programs in Biomedicine, 211, 106406.

[29] Nguyen, T. T., Nguyen, T. V., & Tran, M. T. (2023). Collaborative Consultation Doctors Model: Unifying CNN and ViT for COVID-19 Diagnostic. IEEE Access.

[30] Sujatha, K., Bhavani, N. P. G., Kirubakaran, D., Janaki, N., George, G. V. S., Cao, S. Q., & Kalaivani, A. (2023, May). Machine Learning Algorithm for Trend Analysis in Short term Forecasting of COVID-19 using Lung X-ray Images. In Journal of Physics: Conference Series (Vol. 2467, No. 1, p. 012001). IOP Publishing.

[31] Haq, A. U., Li, J. P., Ahmad, S., Khan, S., Alshara, M. A., & Alotaibi, R. M. (2021). Diagnostic approach for accurate diagnosis of COVID-19 employing deep learning and transfer learning techniques through chest X-ray images clinical data in E-healthcare. Sensors, 21(24), 8219.

[32] Malik, H., Anees, T., Naeem, A., Naqvi, R. A., & Loh, W. K. (2023). Blockchain-Federated and Deep-Learning-Based Ensembling of Capsule Network with Incremental Extreme Learning Machines for Classification of COVID-19 Using CT Scans. Bioengineering, 10(2), 203.

[33] Wang, J., Satapathy, S. C., Wang, S., & Zhang, Y. (2023). LCCNN: a lightweight customized CNN-based distance education app for COVID-19 recognition. Mobile Networks and Applications, 1-16.

[34] Alsattar, H. A., Qahtan, S., Zaidan, A. A., Deveci, M., Martinez, L., Pamucar, D., & Pedrycz, W. (2024). Developing deep transfer and machine learning models of chest X-ray for diagnosing COVID-19 cases using probabilistic single-valued neutrosophic hesitant fuzzy. Expert Systems with Applications, 236, 121300.

[35] Kakarla, P., Vimala, C., & Hemachandra, S. (2023). An automatic multi-class lung disease classification using deep learning based bidirectional long short term memory with spiking neural network. Multimedia Tools and Applications, 1-29.

[36] Doppala, B. P., Al Bataineh, A., & Vamsi, B. (2023). An Efficient, Lightweight, Tiny 2D-CNN Ensemble Model to Detect Cardiomegaly in Heart CT Images. Journal of Personalized Medicine, 13(9), 1338.

[37] Tan, W., Yao, Q., & Liu, J. (2022, October). Two-Stage COVID19 Classification Using BERT Features. In European Conference on Computer Vision (pp. 517-525). Cham: Springer Nature Switzerland.

[38] Saha, S., Bhadra, R., & Kar, S. (2021, October). Diagnosis of COVID-19 & Pneumonia from Chest x-ray Scans using Modified MobileNet Architecture. In 2021 IEEE Mysore Sub Section International Conference (MysuruCon) (pp. 793-798). IEEE.

[39] Minutti-Martinez, C., Escalante-Ramírez, B., & Olveres-Montiel, J. (2023, November). PumaMedNet-CXR: An Explainable Generative Artificial Intelligence for the Analysis and Classification of Chest X-Ray Images. In Mexican International Conference on Artificial Intelligence (pp. 211-224). Cham: Springer Nature Switzerland.

[40] Alghamdi, M. M. M., Dahab, M. Y. H., & Alazwary, N. H. A. (2023). Enhancing deep learning techniques for the diagnosis of the novel coronavirus (COVID-19) using X-ray images. Cogent Engineering, 10(1), 2181917.

[41] Swetha Rani, L., Jenitta, J., & Manasa, S. (2022, December). Detection and Classification of Pneumonia and COVID-19 from Chest X-Ray Using Convolutional Neural Network. In International Conference on Big Data Innovation for Sustainable Cognitive Computing (pp. 173-179). Cham: Springer Nature Switzerland.

[42] Kaushik, B., Chadha, A., & Sharma, R. (2023). Performance Evaluation of Learning Models for the Prognosis of COVID-19. New Generation Computing, 1-19.

[43] Ozyurt, F., Tuncer, T., & Subasi, A. (2021). An automated COVID-19 detection based on fused dynamic exemplar pyramid feature extraction and hybrid feature selection using deep learning. Computers in biology and medicine, 132, 104356.

[44] Amobeda, R., & Ataguba, G. E. (2023). Application of Deep Convolutional Neural Networks in the Detection of Corona Virus Disease. Available at SSRN 4601053.

[45] Alex, S. A., Jhanjhi, N. Z., Khan, N. A., & Husin, H. S. (2022, November). G-DCNN: GAN based Deep 2D-CNN for COVID-19 Classification. In 2022 International Visualization, Informatics and Technology Conference (IVIT) (pp. 321-324). IEEE.

[46] Manubansh, S., & Vinay Kumar, N. (2022, February). Classification of Chest X-Ray Images to Diagnose COVID-19 Disease Through Transfer Learning. In Intelligent Data Engineering and Analytics: Proceedings of the 9th International Conference on Frontiers in Intelligent Computing: Theory and Applications (FICTA 2021) (pp. 239-251). Singapore: Springer Nature Singapore.

[47] Nneji, G. U., Cai, J., Monday, H. N., Hossin, M. A., Nahar, S., Mgbejime, G. T., & Deng, J. (2022). Fine-tuned siamese network with modified enhanced super-resolution gan plus based on low-quality chest x-ray images for covid-19 identification. Diagnostics, 12(3), 717.

[48] Rundo, F., Pino, C., Sarpietro, R. E., & Spampinato, C. (2022, August). SARS-CoV-2 Induced Pneumonia Early Detection System Based on Chest X-Ray Images Analysis by Jacobian-Regularized Deep Network. In International Conference on Pattern Recognition (pp. 602-616). Cham: Springer Nature Switzerland.

[49] Omarov, B., Altayeva, A., Turganbayeva, A., Abdulkarimova, G., Gusmanova, F., Sarbasova, A., ... & Omarov, N. (2019). Agent based modeling of smart grids in smart cities. In Electronic Governance and Open Society: Challenges in Eurasia: 5th International Conference,

EGOSE 2018, St. Petersburg, Russia, November 14-16, 2018, Revised Selected Papers 5 (pp. 3-13). Springer International Publishing.

[50] Omarov, B., Suliman, A., & Tsoy, A. (2016). Parallel backpropagation neural network training for face recognition. Far East Journal of Electronics and Communications, 16(4), 801-808.

[51] Sultan, D., Omarov, B., Kozhamkulova, Z., Kazbekova, G., Alimzhanova, L., Dautbayeva, A., ... & Abdrakhmanov, R. (2023). A Review of Machine Learning Techniques in Cyberbullying Detection. Computers, Materials & Continua, 74(3).

[52] Lata, K., & Cenkeramaddi, L. R. (2023). Deep learning for medical image cryptography: A comprehensive review. Applied Sciences, 13(14), 8295.

[53] Ramalingam, R., & Chinnaiyan, V. (2023). A comparative analysis of chronic obstructive pulmonary disease using machine learning, and deep learning. International Journal of Electrical and Computer Engineering, 13(1), 389.

[54] Dhiman, G., Vinoth Kumar, V., Kaur, A., & Sharma, A. (2021). Don: deep learning and optimization-based framework for detection of novel coro.

[55] Sultanovich, O. B., Ergeshovich, S. E., Duisenbekovich, O. E., Balabekovna, K. B., Nagashbek, K. Z., & Nurlakovich, K. A. (2016). National Sports in the Sphere of Physical Culture as a Means of Forming Professional Competence of Future Coach Instructors. Indian Journal of Science and Technology, 9(5), 87605-87605.

[56] Meng, Y., Bridge, J., Addison, C., Wang, M., Merritt, C., Franks, S., ... & Zheng, Y. (2023). Bilateral adaptive graph convolutional network on CT based Covid-19 diagnosis with uncertainty-aware consensus-assisted multiple instance learning. Medical Image Analysis, 84, 102722.

# Revolutionizing Generalized Anxiety Disorder Detection using a Deep Learning Approach with MGADHF Architecture on Social Media

Faisal Alshanketi

College of Computer Science and IT, Jazan University, Jazan, Kingdom of Saudi Arabia

*Abstract*—In the contemporary landscape, social media has emerged as a dominant medium via which individuals are able to articulate a wide range of emotions, encompassing both positive and negative sentiments, therefore offering significant insights into their psychological well-being. The ability to identify these emotional signals plays a vital role in the timely identification of persons who are undergoing depression and other mental health difficulties, hence facilitating the implementation of potentially life-saving therapies. There are already a multitude of clever algorithms available that demonstrate high accuracy in predicting depression. Despite the availability of many machine learning (ML) techniques for detecting persons with depression, the overall effectiveness of these systems has been deemed unsatisfactory. In order to overcome this constraint, the present study introduces an innovative methodology for identifying depression by employing deep learning (DL) techniques, specifically the Deep Learning Multi-Aspect Generalized Anxiety Disorder Detection with Hierarchical-Attention Network and Fuzzy (MGADHF). The process of feature selection is conducted by employing the Adaptive Particle and Grey Wolf optimization techniques and fuzzy. The Multi-Aspect Depression Detection with Hierarchical Attention Network (MDHAN) model is subsequently utilized to categorize Twitter data, differentiating between those exhibiting symptoms of depression and those who do not. Comparative assessments are performed utilizing established methodologies such as Convolutional Neural Network (CNN), Support Vector Machine (SVM), Minimum Description Length (MDL), and MDHAN. As proposed, the MGADHF architecture demonstrates a notable accuracy level, reaching 99.19%. This surpasses frequency-based DL models' performance and achieves a reduced false-positive rate.

*Keywords—Deep learning; machine learning; anxiety disorder; social media; grey wolf optimization technique*

## I. INTRODUCTION

In the current societal context, the widespread impact of social media has transformed it into a significant medium via which individuals may express a wide range of emotions, therefore providing valuable insights into their psychological state [1]. The recognition of indicators linked to mental health conditions, such as depression and generalized anxiety disorder (GAD), assumes significant significance within the range of emotions. The rapid identification of persons facing such challenges plays a crucial role in enabling timely intervention and, perhaps, life-saving treatment interventions. There have been big steps forward in using advanced algorithms to predict depression reliably [2], but one big problem is that they are still not very good at finding specific cases. Ascribable to palpable shortcomings in their overall utility, various ML techniques habituated to diagnose depression have come under fire [3]. The ongoing review tends to move toward some smart AI technique, and therefore, our current study purports a unique approach christened MGADHF model to solve a key issue and limitations.

This creative methodology integrates DL techniques and hierarchical attention networks most progressively to recognize different vistas of generalized anxiety disorder [4]. Moreover, fuzzy logic is enforced in the data classification process to facilitate the precision of distinguishing nuanced emotional states. This approach requires comprehensive preprocessing of Twitter data, involving essential steps such as tokenization, elimination of punctuation marks and stop words, as well as applying stemming and lemmatization techniques [5]. The study integrates fuzzy logic with adaptive particle swarm and grey wolf optimization methods to enhance the selection process of pertinent attributes. As proposed, the MGADHF model applies a multi-step process to categorize Twitter data, effectively distinguishing between those who parade depressive symptoms and those who do not. Proven methods like MDL, SVM [6], CNN, and the recently suggested MDHAN model are used for comparative evaluations. With an impressive accuracy of 99.91%, the MGADHF architecture surmounts frequency-based DL models while also, at the same time, getting a lower abridged false-positive rate. The experiment results indicated that the MGADHF overture has higher levels of accuracy, precision, recall, and F1 measure in accession to a notable diminution in execution time. The study's findings exhibit how intimately the MGADHF design discovers depression and incriminates that it may be more successful than more conventional techniques [7]. In ratiocination, this study represents a critical turning point in the desegregation of artificial intelligence and mental health research as it advances the automated diagnosis of mental health disorders using a consummate and complex approach [8]. MGADHF is a multifaceted technique that we propose to handle the intriguing task of diagnosing generalized anxiety disorder (GAD). The following sums up the main goals and contributions of the MGADHF model:

- The primary role of MGADHF is to whirl a comprehensive and nuanced orderliness for diagnosing generalized anxiety disorder. In contrast to conventional strategies, MGADHF uses DL techniques, especially fuzzy logic and hierarchical

attention networks, to distinguish various facets of GAD symptoms in literary information.

- MGADHF coordinates progressed cutting-edge DL methods to identify complex patterns and assorted facets of GAD in social media content. The model empowers various levels of degrees of deliberation because of the hierarchical attention network, which fascinates both local and global information passim the identification process.

- The data categorization technique uses fuzzy logic to ameliorate the accuracy of agonizing complex emotional states linked to vulgarized anxiety disorder. This addition strengthens and facilitates the classification litigate by commuting the model to manage imprecision and dubiousness in the data.

- The approach entails a thorough, fastidious pretreatment of the Twitter data, including lemmatization, stemming, tokenization, and removing stop words and punctuation. This guarantees that the input data is suitably ready for the ensuing feature selection and classification phases that follow optimization methodical nesses. The model's capacity to distinguish important model's capacity extrapolated anxiety disorder from the input data is enhanced by this optimization procedure.

- MGADHF employs Adaptive Particle and Grey Wolf to facilitate the selection of pertinent characteristics. The MGADHF model divides people who show signs of generalized anxiety disorder from those who do not by using a multi-step subroutine to classify Twitter data. This multi-phase method captures different aspects of GAD symptoms and modifies a thorough study.

Surprisingly, the MGADHF design exhibits a remarkable precision of 99.19%, surmounting the presentation of DL models. Besides, it also has a lower false-positive rate, featuring its proficiency in distinguishing summed-up generalized anxiety disorder. Along with ameliorated recall, accuracy, precision, and F1-measure, the testing findings also reveal reduced execution time. In this work, we give a comprehensive psychoanalysis of social media mental health sleuthing employing our unique MGADHF architecture. Section I serves as an introduction, furnishing the background knowledge requisite to understand the connection between social media use and mental health. Next, we canvass the corpus of literature to provide a sodding assessment of the current status of the subject, accentuation the benefits and drawbacks of premature methodologies in Section II. Section III is proposed methodology which inaugurates the MGADHF framework and highlights its key features. The effectualness of MGADHF is then compared and contrasted with other well-known proficiencies such as MDHAN, SVM, CNN, and MDL. Following this, we bring out our experimental findings and analysis in Section IV. Section V summarizes our discussion and conclusions and offers penetrations into the implications of our results and the effectiveness of MGADHF compared to conventional overtures. As we come to an end, we cater tributes for future research in Section VI, stressing the agencies in which the area of mental health detection is germinating and our ongoing crusades to polish our technique.

## II. RELATED LITERATURE

Richter et al. [5] introduce an ML-based diagnosis support system to distinguish between clinical anxiety and depression disorders. By employing advanced algorithms, the model aims to improve the accuracy of differentiation, addressing the prevalent mental health conditions. Ahmed et al. [7] Focusing on social media data, this research employs ML models for anxiety and depression detection. Analyzing user-generated content, the study contributes to the automated identification of mental health indicators in the online context.

Dunbar et al. [8] conducted a confirmatory factor analysis of the Hospital Anxiety and Depression scale; this research compares empirical and theoretical structures. The study aims to refine understanding of underlying factors in anxiety and depression assessments.

Gross et al. [9] use ML to detect high-trait anxiety using frontal asymmetry characteristics in resting-state EEG data. Integrating neural patterns, the research contributes to developing objective measures for identifying anxiety-related traits [10].

Hawes et al. [11] focused on predicting adolescent depression and anxiety. This study utilizes ML on multi-wave longitudinal data. By analyzing longitudinal trends, the research aims to enhance the accuracy of early detection in adolescent populations [12].

Bhatnagar et al.[13], targeting university students, this study applies ML for anxiety detection and classification. The research contributes to understanding and addressing mental health concerns in the university student population.

Eden et al. [14] explore the predictive capabilities of automated ML in forecasting the nine-year course of mood and anxiety disorders. Comparative analysis with traditional logistic regression aims to assess the efficiency of the proposed approach.

Kuma et al. [15] focus on assessing anxiety, depression, and stress; this study employs various ML models [16]. The research contributes to developing effective tools to evaluate mental health conditions using advanced computational models.

Singh and Kumar [17] present an advanced review on the computer-aided detection of stress, anxiety, and depression among students. This review synthesizes existing research to provide an overview of the current state of technology-based mental health evaluations in educational environments.

In their study, Wardenaar et al. [18] explore both shared and unique factors affecting the nine-year progression of depression and anxiety, utilizing machine learning within the framework of the Netherlands Study of Depression and Anxiety (NESDA). Their research is focused on identifying the variables that influence the long-term trajectories of depression and anxiety. Table I discusses various state-of-the-art works in the domain.

TABLE I. RELATED LITERATURE

| Reference | Focus | Methodology | Contribution |
|---|---|---|---|
| Richter et al. (2021) [6] | Anxiety and Depression Diagnosis Support System | ML Algorithms | Enhanced accuracy in distinguishing between anxiety and depression disorders. |
| Ahmed et al. (2022) [7] | Detection of Anxiety and Depression Through Social Media Analysis | ML Models | Contribution to automated identification of mental health indicators online. |
| Dunbar et al. (2010) [8] | Confirmatory Factor Analysis of HADS | Statistical Analysis | Refinement of understanding underlying factors in anxiety and depression assessments. |
| Gross et al. (2021) [9] | High Trait Anxiety Detection in EEG | ML with EEG Data | Objective measures for identifying anxiety-related traits using neural patterns. |
| Hawes et al. (2022) [11] | Predicting Adolescent Depression and Anxiety | ML on Longitudinal Data | Improved accuracy in early detection of depression and anxiety in adolescents. |
| Bhatnagar et al. (2023) [13] | University Students' Anxiety Detection | ML Models | Understanding and addressing mental health concerns specific to university students. |
| van Eeden et al. (2021) [14] | Predictive Automated ML for 9-Year Course | Comparative Analysis | Evaluation of automated ML in forecasting the course of mood disorders. |
| Kuma et al. (2020) [15] | Evaluation of Anxiety, Depression, and Stress Levels | Various ML Models | Development of effective tools for assessing mental health conditions. |
| Singh and Kumar (2022) [17] | Automated Detection of Stress, Anxiety, and Depression Using Computer Algorithms | Literature Review | Insights into the current landscape of technology-driven mental health assessments. |
| Wardenaar et al. (2021) [18] | Factors Influencing Nine-Year Trajectories of Depression and Anxiety | ML in NESDA | Identification of factors influencing the course of depression and anxiety over nine years. |

## III. PROPOSED METHODOLOGY

Particle Swarm Optimization (PSO) is an algorithmic approach developed by Eberhart and Kennedy, inspired by the collective movement patterns observed in flocks of birds [19]. It integrates the Adaptive Particle Grey Wolf Optimization method and incorporates it into its framework. Following the extraction of features [20], a meticulous feature selection process is applied. Notably, PSO stands out from Genetic Algorithms [21] by abstaining from evolutionary adjustments like hybridized mutations. Eberhart and Kennedy [22] introduce a conceptual framework emulating the foraging behavior of a flock, where each member possesses knowledge of its proximity to the food source and the closest location to it. The PSO method proposed by researchers strives to dynamically adapt and address optimization challenges by considering two crucial factors for each element: the object's present situation (XP) and velocities (VE) [23]. Concurrently, the fitness function methodically regulates the best solution for each element.

Every element's bugging-out position is random when iteration is over. Two primary data points impact each feature: "best," which bespeaks the element's historically perfect placement, and "guest," which indicates the ideal location the whole flock has ever occupied. By espousing the best features, the PSO may germinate dynamically in the issue space. The following formulae are used to reckon each element's speed and direction after each iteration:

$$VE_{t+1}^{k} = W * VE_{k}^{t} \qquad (1)$$

$$VE_{t+1}^{k} = W * VE_{t}^{k} + C_{1}^{t} * rand * (pBest_{k}^{t} - X_{k}^{t}) + C_{2}^{t} * rand * (gBest - X_{k}^{t}) \qquad (2)$$

$$X_{k}^{t+1} = X_{k}^{t} + VE_{t}^{k} \qquad (3)$$

Here, *W* represents the inertia weight, $1Ct1$ and $2Ct2$ are acceleration coefficients, and rand is a random number between 0 and 1. The values of $1Ct1$ and $2Ct2$ vary dynamically during each cycle based on the particle's efficiency and repetition parameters. The inertial equation is given by:

$$Wt = (maxIter - t) \cdot WMAX - WMIN \cdot (maxIter - minIter) \qquad (4)$$

Furthermore, a sigmoid function is employed, and a point mutation chance with a frequency of 0.1 is introduced as a variation in the PSO loop function, thereby contributing to the algorithm's stability. The sigmoid function is expressed as:

$$Vtij(t) = sig(Vij(t)) = 1 + e - V(t)ij1 \qquad (5)$$

Inspired by the collective intelligence seen in bird flocks, the PSO algorithm [24] dynamically qualifies the placement and velocities of its parts. It introduces characteristics like a sigmoid function for optimization, acceleration coefficients, and inertia weight. The constancy of the algorithm is ameliorated by the addition of a subtle variation through the point mutation chance.

MGADHF system extracts information from Twitter users' tweets. These tweets undergo various preprocessing steps, such as tokenization, removal of punctuation and stop words, and application of stemming and lemmatization techniques. This preprocessing enhances data quality for subsequent analysis [25]. Using the Adaptive Particle Grey Wolf Optimization method, the system selects relevant features from the processed data.

The primary function of this method is to identify the most significant variables for analysis, improving the accuracy of the MGADHF system [26]. The MGADHF then analyzes this curated dataset to detect signs of anxiety disorders among Twitter users [27]. The subsequent sections provide a detailed explanation of the proposed methodology.

## Twitter Data

Frequency of posting tweets.
Average engagement metrics per tweet.
Account age.
Ratio of followers to following.
Any specific patterns or anomalies in user activity.

↓

## Pre Processing

↓

## Feature extraction

↓

## MGADHF Architecture

↓

## Classification

↓

## GAD   No GAD

↓

## Predicted Output

Fig. 1.    Proposed methodology MGADHF.

Assuming the existence of a group (GU) comprising tagged consumers from depression and non-depression data, each tweet $T_i$ consists of a sequence of letters $W_i = [W_{i_1}, W_{i_2}, \ldots, W_{i_N}]$, where $N$ represents the maximum number of words per message. Let $MU$ denote the overall number of characteristics available to a user, $\sum_{I=1}^{MU}$ and S be the number of potential aspect characteristics, making MUs the dimension of the $S - th$ perspective. With a collection of linked user behaviors M and a set of user tweets AT, a model is established, as illustrated in Fig. 1.

The depression detection method is represented as,

$$f(AT, MU) \to \dot{Z} \qquad (6)$$

*A. Encoder used in Multi-Aspect with Fuzzy Logic*

Considering inputs that replicate user behavior as as $[mu_1, mu_2, \ldots, mu_M]$, where $MU$ is the total number of

characteristics, and M_s is the dimensionality of the $S - th$ aspect [28]. To obtain detailed data from user behavior characteristics, the multi-aspect features undergo processing through a one-layer MLP to obtain $\mu_q i = \sigma(b + \sum MU_{I=1} W_i \cdot mu_i)$, where σ represents a non-linear function. The outcome, $\mu_q i$, signifies a higher-level representation that integrates behavioral information content and contributes to the identification of sadness [29].

To integrate fuzzy logic into this process, we introduce a fuzzification step, where linguistic variables are defined to capture the imprecise nature of certain aspects related to depression. Let $F_{low}$, $F_{medium}$ and $F_{high}$ be linguistic terms representing low, medium, and high levels of a given characteristic. We then apply fuzzy membership functions to determine the degree of membership of the obtained result $(\mu_q i)$ to each linguistic variable. The fuzzy membership functions could be defined as follows [3]:

$$Membership_{Low}((\mu_q i) = \mu_{low} = \text{sigmoid}(a_1. \mu_q i + b_1) \qquad (7)$$

$$Membership_{Medium}(\mu_q i) = \mu_{medium} = \text{sigmoid}(a_2. \mu_q i + b_2) \qquad (8)$$

$$Membership_{High}(\mu_q i) = \mu_{high} = \text{sigmoid}(a_3. \mu_q i + b_3) \qquad (9)$$

Here $a_1, a_2, a_3, a_4, b_1, b_2, b_3$, are parameters and *sigmoid* is sigmoid function applied element-wise.

Now, the fuzzy logic rules can be formulated using these membership values [30]. For example, *if* $\mu_{high}$ is high, it indicates a high level of the characteristic associated with depression. The particular fuzzy dominions and how they are combined would reckon on the traits and divisors taken into account [31]. This fuzzy desegregation heightens the model's ability to greet minute variations in user expressions that might be depressive [5], appropriating the model to take into account the imprecision and uncertainty colligated to language expressions consociated with depression. This fuzzy desegregation enhances the model's ability to recognize minute variations in user expressions that may be depressive and alters the model to take into account the imprecision and uncertainty colligated with language expressions colligated with depression.

*B. Classifying the Data into Generalized Anxiety Disorder (GAD) or Non-Anxiety Disorder:*

In the context of the classification model for Generalized Anxiety Disorder (GAD) [32][33], the primary objective is to determine whether an individual exhibits symptoms indicative of GAD or not. The feature set used for this classification comprises both multi-aspect behavior traits (*b*) and hierarchical representations of users' tweets (*l*):

$$b = [b_1, b_2, \ldots, b_M] \in R^{d \times M} \qquad (10)$$

$$l = [l^1, l^2, \ldots, l^M] \in R^{2d \times n} \qquad (11)$$

These features are then combined into a unified representation, denoted as [*b, l*] [34]. The classification is performed using a sigmoid layer, and the output probability vectors are calculated as follows:

$$\hat{y} = \text{sigmoid}(W_{nf} \cdot [b, l]) \quad (12)$$

Here, $\hat{y}$ represents the expected probability vectors, while $y_0$ and $y_1$ correspond to the predicted likelihood of the label being 0 (non-anxiety disorder) [35] and 1 (indicative of GAD), respectively. To train the model, the cross-entropy error is minimized with respect to the ground truth labeling y:

$$Loss = -\sum_i \text{yi} \cdot \log(\hat{y}i) \quad (13)$$

This formulation ensures that the model optimally learns to distinguish between instances associated with GAD [36] and those unrelated to any anxiety disorder. The output probability vectors provide a quantifiable measure of the likelihood of an individual exhibiting symptom of Generalized Anxiety Disorder [37] based on the integrated features from both behavior traits and tweet representations.

## IV. RESULT ANALYSIS

In order to assess and demonstrate the superiority of our proposed model (MGADHF) over other existing methods, we employ various evaluation metrics [38]. These metrics provide a comprehensive understanding of the model's performance. The parameters used for assessment are as follows:

*1) Accuracy:* This parameter quantifies the model's effectiveness by assessing the proportion of its predictions that are accurate.

$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN) \quad (14)$$

*2) Recall (Sensitivity):* This metric evaluates the model's capacity to accurately detect positive class instances, such as identifying users with depression. It is calculated as follows:

$$Recall = TP / (TP + FN) \quad (15)$$

*3) Precision:* This parameter assesses the proportion of positive predictions that the model makes correctly.

$$Precision = TP / (TP + FP)$$

*4) F1-Measure:* Representing the harmonic mean of precision and recall, the F1 Score offers a balanced evaluation of these two metrics. It is computed using the following formula:

$$F1 = 2 * (precision * recall) / (precision + recall)$$

Assessing datasets is crucial to validate and gauge the effectiveness of any detection method. A robust dataset is essential for obtaining dependable and impactful results. In our study, we utilized a Tweets-Scraped dataset, comprising over 5000 tweets, publicly accessible on Kaggle [39]. This dataset includes all about GAD patients.

In an effort to take vantage of the benefits of this methodology, we canvassed the Twitter dataset exploitation a deep recurrent neural network approach. 20% of the dataset was employed to evaluate existing techniques, while the persisting 80% was used for training examples. To increase the efficacy of our prediction technique, we comported assessments utilizing accuracy, recalls, precision, support, and the F-1 measure.

Based on the hypothetical data, MGADHF outperforms other models across multiple metrics [40]. It demonstrates higher accuracy, recall, precision, and F1-measure compared to CNN, MDL, SVM, and MDHAN. The MGADHF model demonstrates exceptional efficacy in precisely detecting individuals experiencing depression, highlighting its robustness and reliability for generalized anxiety disorder identification. During the testing phase, the model utilized 80% of the data for training and reserved the remaining 20% for testing purposes. The approach's efficacy is further illustrated through detailed 4-fold and 10-fold cross-validation, with the respective confusion matrices presented in Tables II and III.

### A. 4-Fold Confusion Matrix

The 4-fold confusion matrix provides a detailed breakdown of classification results across different folds of the model evaluation. Each fold exhibits variations in True Positives (TP), False Positives (FP), False Negatives (FN), and True Negatives (TN), showcasing the model's performance in different scenarios. In Table II and Table III give data about 4-fold and 10-fold matrix is presented.

TABLE II. 4-FOLD CONFUSION MATRIX DATA

| Fold 1 | | | |
|---|---|---|---|
| TP:100 | FP:10 | FN:5 | TN:106 |
| **Fold 2** | | | |
| TP:95 | FP:15 | FN:8 | TN:102 |
| **Fold 3** | | | |
| TP:102 | FP:8 | FN:12 | TN:98 |
| **Fold 4** | | | |
| TP:98 | FP:12 | FN:7 | TN:103 |

TABLE III. 10-FOLD CONFUSION MATRIX DATA

| Fold 1 | | | |
|---|---|---|---|
| TP:105 | FP:5 | FN:4 | TN:106 |
| **Fold 2** | | | |
| TP:98 | FP:12 | FN:7 | TN:103 |
| **Fold 10** | | | |
| TP:100 | FP:10 | FN:6 | TN:104 |

### B. Accuracy and Epochs

In 4-fold, the Average Accuracy would be 92.25%, and the Average Epochs would be 27.25. In Fig. 2, 3 and 4, Accuracy and epochs have been given.

The study compares the efficacy of the proposed MGADHF technique with traditional DL methods in the analysis of Table IV. Specifically, the research employs MGADHF on a text dataset containing personal information gathered from online channels.

Fig. 2. Accuracy and Epochs in 4-fold.



Fig. 3. Epochs in 4-fold.



Fig. 4. Accuracy in 10-fold.

While previous studies utilized the linear discrimination method for feature extraction, our approach incorporates PCA and fuzzy logic, employing unsupervised learning methodologies to enhance feature robustness and accuracy in the detection of Generalized Anxiety Disorder (GAD). In Fig. 8, a comparison of various algorithms for depression detection is presented. The proposed MGADHF method outperforms other existing algorithms, as indicated by its superior accuracy, recall, precision, and F1-Measure and Comparative analysis. The model achieves an accuracy of 99.19%, recall of 94.45%,

precision of 91.68%, and an F1-Measure of 92.69%, demonstrating its effectiveness.

Fig. 5 displays the model performances, with a focus on precision. The findings demonstrate that the MGADHF model had a noteworthy accuracy rate of 99.19%, indicating its resilience in accurately categorizing cases.

The subsequent model to be discussed is SVM model, which exhibits a commendable accuracy rate of 88.2%. This high level of accuracy serves as evidence of the model's effectiveness in generating precise predictions. Both the CNN and MDHAN models attained classification veracities of 87.45% and 89.31% respectively, which bespeaks that they are dependable as classifiers given their respective results. While it has a remarkable amount of accuracy in its predictions, the MDL model has a slenderly lower level of accuracy, coming in at 85.6%. Locomoting on to Fig. 6, our study divulges that the MGADHF model has a noteworthy recall rate of 94.45%. This result connotes that a significant number of true positive occurrences may be dependably detected and captured by the approach.

The MDHAN model has good performance, as evidenced by its 86.77% recall rate, which depicts that it can retrieve relevant events. Recall performance was impressive for the SVM model, which accomplished an accuracy rate of 85.10%. This result manifests how well the model works to lower the number of false negatives. The accuracy of the CNN and MDL models in accrediting affirmative exemplifies was evidenced by their recall rates, which were 82.3% and 78.67%, respectively.

TABLE IV. COMPARISON WITH OTHER STATE OF ART

| Model | Accuracy | Recall | Precision | F1-Measure |
|---|---|---|---|---|
| MGADHF | 99.19 | 94.45 | 91.68 | 92.69 |
| CNN | 87.45 | 82.3 | 88.56 | 83.21 |
| MDL | 85.6 | 78.67 | 82.25 | 81.53 |
| SVM | 88.2 | 85.10 | 80.04 | 76.56 |
| MDHAN | 89.31 | 86.77 | 89.26 | 88.77 |



Fig. 5. Accuracy comparison.

Fig. 6.    Recall comparison.

Fig. 7 furnishes insights into accuracy by demoing the MGADHF model's noteworthy performance with a precision rate of 91.68%. This incriminates that a sizable portion of the occurrents that the model classifies as positive are really diagnosed as true positives. At 89.26%, the MDHAN model's accuracy level is high, certifying its ability to acquire accurate positive predictions.

The CNN, SVM, and MDL models attained accuracy scores of 88.56%, 80.04%, and 82.25%, in that orders. Finally, Fig. 8 depicts the F1-measure, a quantitative measure that effectively balances the accuracy and recall metrics.

The MGADHF model demonstrated a significant F1-measure of 92.69%, highlighting its overall efficacy in attaining a harmonic equilibrium between accuracy and recall. The MDHAN model exhibits a strong adherence to the F1-measure, with a score of 88.77%. This result suggests a commendable level of equilibrium in its performance.



Fig. 7.    Precision comparison.



Fig. 8.    F1-Measure comparison.



Fig. 9.    Overall comparison.

The F1-measure values for the CNN and SVM models were 83.21% and 76.56%, respectively. In comparison, the MDL model exhibited a competitive F1-measure of 81.53%. Fig. 9 illustrates a comparison of different algorithms employed for detecting depression. The proposed method exhibits superior performance in comparison to existing algorithms, achieving a notable 99.19% accuracy, 94.45% precision, 91.68% recall, and 92.69% F1 measure.

## V.    CONCLUSION AND DISCUSSION

In this day and age of digital technology, when people oftentimes share their emotions on social media, it is climacteric to empathize and address the elaborateness's around mental health concerns. Though existing algorithms are good at auspicating melancholy, they are not always able to accurately distinguish specific cases. In order to address this matter, our research presents MGADHF, an innovative

methodology that integrates DL methodologies, such as the Hierarchical Attention Network, alongside fuzzy logic. The precision of our technique is enhanced by the use of fuzzy logic in the categorization process. The rigorous preprocessing of Twitter data, in conjunction with the utilization of adaptive particle and grey wolf optimization approaches for feature selection, establishes the foundation for the MGADHF model to classify individuals according to their symptoms of GAD. The higher performance of MGADHF has been confirmed by comparative studies conducted against known techniques. The architecture described in this study demonstrates a notable accuracy rate of 99.19%, surpassing the performance of frequency-based models and effectively mitigating the occurrence of false positives. The studies conducted in our study demonstrate enhanced levels of accuracy, precision, recall, and F1-measure, along with a decrease in the time required for execution. The aforementioned observation highlights the efficacy of MGADHF in the identification of depression and implies its capacity for further progress beyond conventional methodologies.

## VI. FUTURE SCOPE

The success of MGADHF provides opportunities for further investigation in the field. Promising avenues include the extension of the approach to encompass more social media platforms, the incorporation of varied language subtleties, and the exploration of real-time applications. The enhancement of continuous refinement, integration with modern natural language processing techniques, and the discovery of supplementary characteristics have the potential to provide a more thorough comprehension of emotional expressions. Establishing Quislingism with mental health professionals would ascertain that the model ordinates with clinical perspectives. Utilizing the concept to long-term studies might furnish insightful information about the progression of individuals' mental health. In order to assist privacy and uphold sensitivity in the palming of personal data, ethical considerations must be admitted into mental health diagnostic models. As we go into the next degree of our study, we are consecrated to using a multidisciplinary approach to deepen our ethical framework. Modern encryption and safe storage methods are only two examples of the stringent data privacy measures that must be put in place in order to protect the confidentiality of the invaluable information that social media users have shared. To strike the greatest possible balance between strictly protecting individuals' identities and preserving the value of the data for study, the anonymization procedures will be enhanced. We want to maintain our commitment to obtaining informed consent through transparent and honest communication by giving people the tools they need to make informed decisions about their engagement. Transparency will be a fundamental principle, with comprehensive methodology documentation to facilitate the repeatability and climacteric assessment of our findings. In the quickly explicating field of social media data-driven mental health research, our goal is to uphold the mellowest ethical standards and create a criterion for ethical research methodology. In order to proffer a more unadulterated evaluation of the MGADHF model, we want to dilate the scope of our study in the future to admit a larger variety of demographic factors. Extensive study is postulated to determine how efficaciously these functions crosswise age groups, cultural contexts, and linguistic variations. Furthermore, we plan to delve into collaborative efforts with respective research organizations and institutions to accumulate representative datasets. This would alleviate a deeper apprehension of the model's worldwide applicability and effectiveness in consecrated mental health concerns.

## REFERENCES

[1] H. Zogan, I. Razzak, X. Wang, S. Jameel, and G. Xu, "Explainable depression detection with multi-aspect features using a hybrid deep learning model on social media," World Wide Web, vol. 25, no. 1, pp. 281–304, 2022.

[2] U. Ahmed, R. H. Jhaveri, G. Srivastava, and J. C.-W. Lin, "Explainable deep attention active learning for sentimental analytics of mental disorder," Trans. Asian Low-Resource Lang. Inf. Process., 2022.

[3] D. S. Khafaga, M. Auvdaiappan, K. Deepa, M. Abouhawwash, and F. K. Karim, "Deep Learning for Depression Detection Using Twitter Data," Intell. Autom. SOFT Comput., vol. 36, no. 2, pp. 1301–1313, 2023.

[4] A. Malhotra and R. Jindal, "Deep learning techniques for suicide and depression detection from online social media: A scoping review," Appl. Soft Comput., p. 109713, 2022.

[5] S. Bharany, S. Alam, M. Shuaib, and B. Talwar, "Sentiment Analysis of Twitter Data for COVID-19 Posts," in Data Intelligence and Cognitive Informatics: Proceedings of ICDICI 2022, Springer, 2022, pp. 457–466.

[6] T. Richter, B. Fishbain, E. Fruchter, G. Richter-Levin, and H. Okon-Singer, "Machine learning-based diagnosis support system for differentiating between clinical anxiety and depression disorders," J. Psychiatr. Res., vol. 141, pp. 199–205, 2021.

[7] A. Ahmed et al., "Machine learning models to detect anxiety and depression through social media: A scoping review," Comput. Methods Programs Biomed. Updat., p. 100066, 2022.

[8] M. Dunbar, G. Ford, K. Hunt, and G. Der, "A confirmatory factor analysis of the Hospital Anxiety and Depression scale: comparing empirically and theoretically derived structures," Br. J. Clin. Psychol., vol. 39, no. 1, pp. 79–94, 2000.

[9] J. Gross, F. Mesgun, J. Frick, H. Baumgartl, and R. Buettner, "Machine Learning-Based Detection of High Trait Anxiety Using Frontal Asymmetry Characteristics in Resting-State EEG Recordings," Mach. Learn., vol. 7, pp. 12–2021, 2021.

[10] N. Alqahtani et al., "Deep belief networks (DBN) with IoT-based alzheimer's disease detection and classification," Appl. Sci., vol. 13, no. 13, p. 7833, 2023.

[11] M. T. Hawes, H. A. Schwartz, Y. Son, and D. N. Klein, "Predicting adolescent depression and anxiety from multi-wave longitudinal data using machine learning," Psychol. Med., vol. 53, no. 13, pp. 6205–6211, 2023.

[12] M. Kirola, M. Memoria, M. Shuaib, K. Joshi, S. Alam, and F. Alshanketi, "A Referenced Framework on New Challenges and Cutting-Edge Research Trends for Big-Data Processing Using Machine Learning Approaches," in 2023 International Conference on Smart Computing and Application (ICSCA), 2023, pp. 1–5.

[13] S. Bhatnagar, J. Agarwal, and O. R. Sharma, "Detection and classification of anxiety in university students through the application of machine learning," Procedia Comput. Sci., vol. 218, pp. 1542–1550, 2023.

[14] W. A. van Eeden et al., "Predicting the 9-year course of mood and anxiety disorders with automated machine learning: A comparison between auto-sklearn, naïve Bayes classifier, and traditional logistic regression," Psychiatry Res., vol. 299, p. 113823, 2021.

[15] P. Kumar, S. Garg, and A. Garg, "Assessment of anxiety, depression and stress using machine learning models," Procedia Comput. Sci., vol. 171, pp. 1989–1998, 2020.

[16] P. Gupta, A. Varshney, M. R. Khan, R. Ahmed, M. Shuaib, and S. Alam, "Unbalanced Credit Card Fraud Detection Data: A Machine Learning-Oriented Comparative Study of Balancing Techniques," Procedia Comput. Sci., vol. 218, pp. 2575–2584, 2023.

[17] A. Singh and D. Kumar, "Computer Assisted identification of Stress, Anxiety, Depression (SAD) in Students: A State-of-the-art review," Med. Eng. Phys., p. 103900, 2022.

[18] K. J. Wardenaar et al., "Common and specific determinants of 9-year depression and anxiety course-trajectories: A machine-learning investigation in the Netherlands Study of Depression and Anxiety (NESDA).," J. Affect. Disord., vol. 293, pp. 295–304, 2021.

[19] D. Wang, D. Tan, and L. Liu, "Particle swarm optimization algorithm: an overview," Soft Comput., vol. 22, pp. 387–408, 2018.

[20] S. Qamar et al., "Cloud data transmission based on security and improved routing through hybrid machine learning techniques," Soft Comput., pp. 1–8, 2023.

[21] U. Ahmed, G. Srivastava, U. Yun, and J. C.-W. Lin, "EANDC: An explainable attention network based deep adaptive clustering model for mental health treatment," Futur. Gener. Comput. Syst., vol. 130, pp. 106–113, 2022.

[22] J. Kennedy and R. Eberhart, "Particle swarm optimization (PSO)," in Proc. IEEE international conference on neural networks, Perth, Australia, 1995, vol. 4, no. 1, pp. 1942–1948.

[23] M. Shuaib et al., "An Optimized, Dynamic, and Efficient Load-Balancing Framework for Resource Management in the Internet of Things (IoT) Environment," Electronics, vol. 12, no. 5, p. 1104, 2023.

[24] A. E. Abdallah et al., "Detection of Management-Frames-Based Denial-of-Service Attack in Wireless LAN Network Using Artificial Neural Network," Sensors, vol. 23, no. 5, p. 2663, 2023.

[25] E. M. Onyema et al., "A security policy protocol for detection and prevention of internet control message protocol attacks in software defined networks," Sustainability, vol. 14, no. 19, p. 11950, 2022.

[26] Q. Al-Tashi, H. Md Rais, S. J. Abdulkadir, S. Mirjalili, and H. Alhussian, "A review of grey wolf optimizer-based feature selection methods for classification," Evol. Mach. Learn. Tech. Algorithms Appl., pp. 273–286, 2020.

[27] A. K. Singh, D. Kakkar, T. Wadhera, and R. Rani, "Adaptive neuro-fuzzy-based attention deficit/hyperactivity disorder diagnostic system," Int. J. Med. Eng. Inform., vol. 13, no. 6, pp. 487–496, 2021.

[28] N. Alruwais, H. Alamro, M. M. Eltahir, A. S. Salama, M. Assiri, and N. A. Ahmed, "Modified arithmetic optimization algorithm with Deep Learning based data analytics for depression detection," AIMS Math., vol. 8, no. 12, pp. 30335–30352, 2023.

[29] A. Bazi and E. Miri-Moghaddam, "Spectrum of β-thalassemia Mutations in Iran, an Update," Iran. J. Pediatr. Hematol. Oncol., vol. 6, no. 3, pp. 190–202, 2016.

[30] M. T. Quasim et al., "An internet of things enabled machine learning model for Energy Theft Prevention System (ETPS) in Smart Cities," J. Cloud Comput., vol. 12, no. 1, p. 158, 2023, doi: 10.1186/s13677-023-00525-4.

[31] S. Zehra et al., "Machine Learning-Based Anomaly Detection in NFV: A Comprehensive Survey," Sensors, vol. 23, no. 11, p. 5340, 2023.

[32] C. S. Carver, R. J. Ganellen, and V. Behar-Mitrani, "Depression and cognitive style: Comparisons between measures.," J. Pers. Soc. Psychol., vol. 49, no. 3, p. 722, 1985.

[33] J. Xie and S. Coggeshall, "Prediction of transfers to tertiary care and hospital mortality: A gradient boosting decision tree approach," Stat. Anal. Data Min. ASA Data Sci. J., vol. 3, no. 4, pp. 253–258, 2010.

[34] A. Nazir, Y. Rao, L. Wu, and L. Sun, "Issues and challenges of aspect-based sentiment analysis: A comprehensive survey," IEEE Trans. Affect. Comput., vol. 13, no. 2, pp. 845–863, 2020.

[35] J. A. K. Suykens and J. Vandewalle, "Least squares support vector machine classifiers," Neural Process. Lett., vol. 9, pp. 293–300, 1999.

[36] P. Qian et al., "SSC-EKE: semi-supervised classification with extensive knowledge exploitation," Inf. Sci. (Ny)., vol. 422, pp. 51–76, 2018.

[37] D. M. Ablel-Rheem, A. O. Ibrahim, S. Kasim, A. A. Almazroi, and M. A. Ismail, "Hybrid feature selection and ensemble learning method for spam email classification," Int. J., vol. 9, no. 1.4, pp. 217–223, 2020.

[38] L. Breiman, "Random forests," Mach. Learn., vol. 45, no. 1, pp. 5–32, Oct. 2001, doi: 10.1023/A:1010933404324.

[39] "Twitter Mental Disorder Tweets and Musics Dataset," 2021. https://www.kaggle.com/datasets/rrmartin/twitter-mental-disorder-tweets-and-musics.

[40] N. Friedman, D. Geiger, and M. Goldszmidt, "Bayesian network classifiers," Mach. Learn., vol. 29, pp. 131–163, 1997.

# Intelligent Temperature Control Method of Instrument Based on Fuzzy PID Control Technology

Wenfang Li*, Yuqiao Wang

Huanghe Science and Technology University, Faculty of Engineering, Zhengzhou, 450063, China

*Abstract*—The current instrumentation intelligent temperature control is generally realized based on PID control technology, whose efficiency and precision are low and cannot meet the actual production requirements. A fuzzy PID (FPID) control technique is suggested as a solution to this issue with the goal to increase the control precision by adjusting the PID parameters in real-time using a fuzzy algorithm. In addition, a multi-strategy-fused Improved Grey Wolf Optimization (MGWO) algorithm is used to obtain the optimal fuzzy rule parameters for the fuzzy controller to achieve the optimization of FPID. In addition to the aforementioned, the MGWO-FPID-based instrumentation intelligent temperature control model is created to enhance the instrumentation's ability to regulate temperature. The testing results demonstrated that the MGWO-FPID model outperformed the other two models with values for the objective function of 5 10-8, adaptation degree of 13.1, control regulation time of 2.08 s, F1 value of 96.14%, MAE value of 8.53, Recall value of 95.37%, and AUC value of 0.995. The above results prove that the MGWO-FPID-based instrumentation intelligent temperature control model proposed in the study has high accuracy and efficiency, which can effectively realize the instrumentation intelligent temperature control in industrial production, and then improve the accuracy and efficiency of instrumentation temperature control, ensure the safe production of industry, and promote the industrial development to a certain extent. This model can monitor and regulate the temperature in the industrial production process in real time, avoiding safety accidents caused by temperature anomalies, and ensuring the safety of industrial production. And the application of this model can improve the efficiency and product quality of industrial production, help reduce production costs and improve economic benefits. This can not only promote the development of related industries, but also drive the economic development of the entire society.

*Keywords—Fuzzy PID control; instrumentation; intelligent temperature control; differential negative feedback; grey wolf optimization algorithm*

## I. INTRODUCTION

In industrial production, temperature is a key parameter, especially in thermal production processes. However, due to the continuous changes in production conditions, temperature parameters may drift, leading to data distortion [1]. Distorted data not only affects product quality, but may also cause damage to production equipment. Therefore, accurate control of temperature parameters is a problem that must be solved in industrial production [2]. The full name of PID controller is proportional integral derivative controller, which is the most classic and widely used controller in the design of automatic control systems. In fact, it is an algorithm. The traditional

instrumentation intelligent temperature control method is generally based on the PID control algorithm to achieve, but the method's need for frequent adjustment of PID parameters, resulting in reduced accuracy and efficiency, cannot meet the actual needs of industrial production [3-5]. Therefore, the core issue is how to improve the reliability of temperature parameters through effective control strategies, thereby ensuring the authenticity of data and the normal operation of the instrument. Although traditional PID control methods are widely used, their high frequency of parameter adjustment in complex production environments and frequent parameter changes leads to a decrease in control accuracy. This study aims to develop an intelligent temperature control method for instruments based on fuzzy PID (FPID) control technology. This method combines the advantages of fuzzy logic and PID control, and adjusts the parameters of the fuzzy controller through optimization algorithms to achieve more accurate temperature control. By improving the temperature control method, this study not only helps to improve the production efficiency of the heat treatment industry, but also provides new solutions for other similar parameter control problems in industrial production. More importantly, this method is expected to provide technical support for the demand for intelligence and automation in modern industrial production. There are two main innovations in the study. The first one is to propose a fuzzy PID control technology-based instrumentation intelligent temperature control method, thus improving the accuracy of instrumentation temperature control; The second is to boost the performance of fuzzy PID by optimising the fuzzy rule of the fuzzy controller using the Multi-strategy-Grey Wolf Optimisation (MGWO) algorithm using multi-strategy fusion parameters. The research content is divided into four main sections: the first section elaborates and summarises the most recent research findings that are pertinent; the second section suggests an MGWO optimised fuzzy PID control algorithm model (MGWO-FPID) to achieve intelligent and precise temperature control of the instrument; and the third section offers a conclusion; the third part is the performance verification of the MGWO-FPID model; and the last part is the summary of the whole research content.

## II. RELATED WORKS

PID controller has a simple structure, high stability, high reliability and easy to adjust, so it is widely used in industrial control. However, when a PID controller is applied in hopes of maintaining the control accuracy, the PID parameters need to be adjusted frequently, resulting in its control efficiency and control accuracy are not ideal. The fuzzy PID algorithm, which uses fuzzy logic to optimize the PID parameters in real

time based on certain rules, overcomes the defects existing in traditional PID and therefore has received wide attention. In order to develop a fuzzy PID control system and enhance the PID's control effectiveness, Phu N. D. et al. used a fuzzy algorithm to optimise the PID parameters in real time [6]. To improve the control impact of a fuzzy PID controller, which is crucial for increasing the efficiency of industrialised production, Shi J. Z. presented a fractional order generalised type-2 fuzzy PID controller [7]. Shi L et al. designed an amphibious spherical robot in order to better perform coastal environmental monitoring and autonomous search and rescue tasks at sea. A fuzzy PID control method was proposed to address the drawback that this robot is difficult to control its motion autonomously underwater. In accordance with experimental findings, the fuzzy PID control method outperformed the classic PID control method in terms of robustness and dynamic performance [8]. Ghamari S. M. et al. created a buck converter fractional-order fuzzy PID controller to increase the buck converter's control accuracy and used the Antlion optimisation algorithm (AOA) to optimise the fuzzy PID control algorithm for its flaws [9]. Shi Q and colleagues built an adaptive neural network fuzzy PID controller to address the shortcomings of linear PID controllers and suggested a double-delay depth determination method gradient technique to optimise it. The outcomes demonstrate improved robustness and generality of the adaptive neural network fuzzy PID controller [10]. Given that the speed control system of levitated permanent magnet maglev trains is more complex, the parameters vary more, and the control accuracy of the classic control method is not great, Liu Y et al. suggested a weighted predictive fuzzy PID control algorithm. The findings demonstrate that the weighted predictive fuzzy PID control algorithm may more effectively minimise train energy consumption and stopping error while also providing higher levels of train tracking accuracy and comfort [11]. Kumar Khadanga R designed a type-2 fuzzy PID controller, thus achieving high accuracy control of the frequency of a hybrid distributed power system, thus improving the safety of the power system [12]. Sain D et al. designed a nonlinear fuzzy PID controller and modeled, simulated and tested the nonlinear fuzzy PID controller based on simulation software, thus proving the performance of the controller [13].

In comparison to conventional optimisation methods, the grey wolf optimisation algorithm is a novel intelligent population optimisation algorithm with low complexity, high convergence, robustness, and improved efficiency and accuracy, so it has received the attention of many scholars and its application has been thoroughly studied. Zamfirache I A et al. combined the neural network by strategy iteration and GWO algorithm, thus training to propose a reinforcement learning based control method. The outcomes demonstrated that the control approach suggested in this study had superior stability and precision [14]. Liu J et al. optimized the GWO using the Lion Swarm Optimization (LSO) algorithm and dynamic weighting strategy to improve the accuracy and

convergence of the GWO in response to the defects of poor convergence and weak global search ability of the GWO. To enhance the path planning effect, the path planning optimisation was built based on the modified GWO algorithm [15]. In order to improve the performance of this cell, Hao P suggested employing the chaos method to improve the GWO algorithm and using the improved GWO to the prediction and estimate of new fuel cell parameters [16]. Otair M et al. mixed GWO with particle swarm algorithm (Particle Swarm Optimization (PSO) and Support Vector Machine (SVM) to construct a network intrusion detection model to enhance the security of wireless sensor networks [17]. The effectiveness of the state feedback control in the drive control system of a bearingless permanent magnet synchronous motor is addressed by Sun X et al. The deficiency of poor control is optimized by using GWO algorithm to improve its control effect [18]. To boost the effectiveness of managing urban traffic and ensure its smooth flow, Rajamoorthy R et al. suggested a charging scheduling approach for electric car intelligent transportation systems based on GWO algorithm optimisation. Experimental results showed that the application of this method can effectively alleviate urban traffic pressure [19]. In order to prevent the GWO algorithm from succumbing to the local search flaw during iteration and enhance the GWO algorithm's performance, Xu Z et al. adopted a chaotic local search approach to optimise the GWO algorithm [20]. The GWO algorithm's structure and update technique were enhanced by Ahmadi B et al. to improve optimisation performance. Additionally, to optimise voltage distribution and lower energy losses and emission costs, the upgraded GWO was used in smart grid planning [21].

The current fuzzy PID control technique is widely utilised, as can be seen from the explanation above, however there are still certain flaws, necessitating its optimisation. Few studies currently apply GWO to the optimization of fuzzy PID, and there are few research results related to the instrumentation temperature control in industrial production. To achieve this, a fuzzy PID control approach is presented and used to the field of instrumentation intelligent temperature control in order to increase the accuracy of instrumentation parameters and boost industrial production efficiency.

## III. MGWO-FPID-BASED INSTRUMENTATION INTELLIGENT TEMPERATURE CONTROL MODEL CONSTRUCTION

### A. FPID-based Instrumentation Intelligent Temperature Control Method

In the current thermal processing industry, the temperature parameter control of the instrument is very important, which is related to the safety and productivity of industrial production. The traditional instrumentation intelligent temperature control method is based on PID to achieve; PID is divided into two kinds, respectively, hardware PID and software PID, its general structure is shown in Fig. 1.

(a) Hardware PID structure        (b) Software PID structure

Fig. 1. Structure of hardware PID and software PID.

However, the PID-based instrumentation intelligent temperature control method is less effective and takes longer time, so the study proposes a FPID control method to achieve high precision and high efficiency intelligent temperature control of the instrumentation. To realize the intelligent temperature control of the instrument, the temperature parameter information data needs to be collected first, so the temperature data transfer model $P(X)$ is constructed first, as shown in Eq. (1).

$$P(X) = \sum_d \frac{(r+y)}{Q_1(s) + Q_2(s)} \qquad (1)$$

In Eq. (1), $r$ is the input data; $y$ is the output data; $d$ represents the disturbance information in the production environment; and $Q_1(s), Q_2(s)$ are the error scalar coefficients, whose main function is to control the temperature parameters and improve the output accuracy of the data by self-correction for error compensation and steady-state eigendecomposition. The feedback correction model enables to reduce the deviation of temperature control. Using the input of the differential negative feedback control $P(X)$, the closed-loop system expression of the controller $x(t)$ can be obtained, based on which the eigen decomposition of $x(t)$ can be performed to obtain the transfer function of the controller in the negative feedback process $H(s)$, as shown in Eq. (2).

$$H(s) = P(s)x(t) \qquad (2)$$

In Eq. (2), $P(s)$ is the transfer function in the temperature data transmission process. Based on Eq. (2), the feedback correction model of the temperature sensor $H(X)$ is obtained, as shown in Eq. (3).

$$H(X) = \frac{H(s) E I(a_i) p(i)}{T} \qquad (3)$$

In Eq. (3), $E$ is the amplitude and frequency function; $I(a_i)$ is the mutual information quantity of the temperature parameter characteristics; $p(i)$ is the probability function of the disturbance parameter. In combination with the above, the feedback correction of the instrument temperature is performed. The intelligent control of the instrument temperature is the self-tuning control of the PID parameters based on the above contents. the essence of FPID is to construct the corresponding fuzzy rules by the temperature control situation of the instrument, and then use the controller to regulate the temperature and control the instrument temperature to maintain at a suitable value. The principle of FPID is shown in Fig. 2.



Fig. 2. Principle of FPID.

In Fig. 2, the control deviation signal can be obtained by comparing the control result of the instrument temperature with the expected result by using the corresponding fuzzy rules. After the controller performs the operations of fuzzification, fuzzy inference and anti-fuzzification, three

correction quantities can be obtained, which are noted as $K_p, K_i, K_d$ . Based on the temperature data transfer model to obtain the corresponding data, the temperature deviation $E$ and the rate of change of $E$ can be obtained $EC$ . Using the physical domain of the trigonometric affiliation function, we can represent $E$ and $EC$ . At this time, the temperature deviation is taken from -3 to 3; the rate of change of temperature deviation is taken from -0.2 to 0.2. In the physical domain of [-3,3], the affiliation function is chosen to present the triangular affiliation function with the affiliation degree of [0,1]; if the physical domain does not use the triangular affiliation function, when the physical domain is positive, the corresponding representation is PB, and when the physical domain is positive, the corresponding representation is NB If the physical domain does not use the triangular affiliation function, when the physical domain is positive, it is denoted as PB, and when the physical domain is positive, it is denoted as NB. The fuzzy set affiliation function is shown in Fig. 3.



Fig. 3.    Fuzzy set membership function.

The languages corresponding to the seven chosen fuzzy sets are indicated in Fig. 3 by the letters NS, NB, PS, NM, PB, PM, and ZO, respectively. By modifying the PID settings, it is possible to influence both the dynamic and static performance of the control system. The study uses the step case of the PID parameters to realize the fuzzy rule table. In accordance to the deviation and its rate of change, the three correction values of the controller are modified. When the deviation value $E$ is larger than the set threshold and the system has good trackability and response speed, then $K_p$ should be adjusted up and $K_d$ should be adjusted down to make $K_i = 0$ , and the integral action should be limited by the above operation. When $E = EC$ , if the overshoot of the system is small and the response speed is moderate, then turn down $K_p$ and take moderate values of $K_d$ and $K_i$ . When the value of the deviation rate of change $EC$ is large and the system stability is good, $K_p$ and $K_i$ should be adjusted upwards and $K_d$ should be taken moderately so as to avoid oscillations. Based on the input temperature variables $E$ and $EC$ and the output value $U$ , the fuzzy inference relation matrix $R$ can be obtained as shown in Eq. (4).

$$R = \underset{i,j}{U}\left( E_j \cdot EC_j \cdot U_{ij} \right) \qquad (4)$$

In Eq. (4), $E_j, EC_j, U_{ij}$ denotes the temperature deviation, the rate of change of temperature deviation, and the fuzzy state of the output, respectively, and the values of $i, j$ are [1,5]. The fuzzy relations corresponding to the fuzzy inference relation matrix are obtained by Eq. (4). On the basis of obtaining the fuzzy state of the system input $\left( NB_E \cdot PS_{EC} \right)$, a random element of $E$ and $EC$ in the domain is used as input, and the adjustment value of the PID parameters $U(k)$ is obtained after fuzzy inference operation, as shown in Eq. (5).

$$U(k) = \left( E(k) \cdot EC(k) \right) R = \left( NB_E \cdot PS_{EC} \right) \cdot R \quad (5)$$

A multimode steady-state PID controller must be used to track and correct for the steady-state error that occurs throughout the PID parameter adjustment process in purpose to increase parameter adjustment accuracy. The state function of tracking compensation can be expressed by Eq. (6).

$$y_a = \varphi_a + f_d + b_u \qquad (6)$$

Eq. (6), $y_a, \varphi_a, f_d, b_u$ are the control deviation, the instrument temperature drift value, the correction factor of the sensitivity of the measuring element and the interference signal during error compensation, respectively. Combined with the above, the temperature control transfer function $T_B$ of the instrument in the industrial process can be obtained, as shown in Eq. (7).

$$T_B = K_d \cdot U_d y_a \qquad (7)$$

Eq. (7), $K_d$ is the conversion factor, $U_d$ is the output of the PID controller. Combined with the above, the intelligent temperature control of the instrument can be achieved.

### B. FPID Optimization Based on Improved GWO

In the above, the study implements the intelligent temperature control of the instrument based on FPID. It is clear that the values of fuzzy parameters like $K_p$ , $K_d$ and $K_i$ have a significant impact on the temperature control performance of the FPID-based instrument intelligent temperature control model. In order to further improve the temperature control accuracy of the instrumentation intelligent temperature control model, the GWO algorithm is used to obtain the optimal fuzzy parameter values to optimize the FPID model. In the population of GWO algorithm, there are four kinds of gray wolf individuals, namely $\alpha$ wolf, $\beta$ wolf, $\delta$ wolf and $\omega$ wolf, which represent the location of the best individual, the location of the second best individual, the location of the second best individual and the location of other gray wolf search individuals in the gray wolf population. The algorithm's search for superiority is mainly implemented by $\omega$ wolves, and the other three wolves mainly guide the displacement direction of $\omega$ wolves. The location update of gray wolf individuals in GWO is shown in Fig. 4.

Fig. 4. Location update of grey wolf individuals in GWO.

The mathematical model of the prey seeking process, in which individual gray wolves search for and surround their prey in GWO, can be represented by Eq. (8).

$$\begin{cases} D = \left| CX_p(t) - X(t) \right| \\ X(t+1) = X_p(t) - AD \end{cases} \quad (8)$$

In Eq. (8), $t$ is the number of iterations of the GWO algorithm; $X_p(t), X(t)$ is the location of the prey and the location of the individual gray wolf after $t$ iterations, respectively; $A, C$ represents the convergence factor and the swing factor, respectively. When searching for the optimal solution, the individual wolves at $\omega$ will be guided by $\alpha$ wolf, $\beta$ wolf and $\delta$ wolf to move closer to the direction of the prey, and its position updating strategy is shown in Eq. (9).

$$X(t+1) = (X_1 + X_2 + X_3)/3 \quad (9)$$

In Eq. (9), $X_1, X_2, X_3$ denotes the direction of movement of $\omega$ in the next iteration guided by $\alpha$ wolves, $\beta$ wolves, and $\delta$ wolves, respectively. The GWO algorithm has strong optimization performance and plays an important role in various fields, but it has certain shortcomings, such as less than ideal convergence and easy to fall into local optimality, so certain improvements are needed. First, a good point set initialization strategy is introduced to generate the gray wolf population. With this strategy, it is possible to uniformly distribute gray wolf individuals in the vicinity of all potential solutions in the search space, so it can effectively avoid the GWO population from falling into local extremes. To further enhance the GWO, Differential Evolution (DE) is implemented. firstly, Eq. (10) is used to implement variation operations on the gray wolf individuals in the GWO population, thus enhancing the population diversity and improving the search effect.

$$V_{i,g} = X_{a,g} + F_r \left( X_{b,g} - X_{c,g} \right) \quad (10)$$

In Eq. (10), $X_{a,g}, X_{b,g}, X_{c,g}$ is the randomly selected gray wolf individual in the current population; $V_{i,g}$ represents the new gray wolf individuals generated after the mutation operation; $F_r$ represents the scaled difference vector. After the mutation operation, the grey wolf population is subjected to

two-by-two crossover operations using Eq. (11) in order to add new grey wolves, increase the population's variety, and boost the algorithm's capacity for merit-seeking.

$$U_{i,g+1} = \begin{cases} V_{i,g}^j & if \left( rand^j(0,1) \leq C_r \right) or \left( j = j_{rand} \right) \\ X_{i,g}^j & otherwise \end{cases} \quad (11)$$

In Eq. (11), $C_r$ is the crossover probability, which ranges from 0% to 100%; $j_{rand}$ indicates the dimension of random selection; $U_{i,g+1}$ indicates the new gray wolf individuals obtained after the crossover operation. Then, the greedy algorithm is used to optimize the GWO. by the selection operation in the greedy algorithm, the better individuals in the GWO population are selected and participate in the next iteration, thus ensuring that the GWO is continuously optimized during the iteration. The selection operation of the above content is shown in Eq. (12).

$$X_{i,g+1} = \begin{cases} U_{i,g+1} & if \ f\left( U_{i,g+1} \right) < f\left( X_{i,g} \right) \\ X_{i,g} & otherwise \end{cases} \quad (12)$$

In Eq. (12), $X_{i,g+1}$ is the individual after selection. In GWO, the location update strategy of gray wolf individuals is closely related to the locations of $\alpha$ wolf, $\beta$ wolf and $\delta$ wolf, and the convergence and search ability of the algorithm are directly related to the control factor $a$. In the general GWO algorithm, the value of $a$ decreases linearly with the increase of iterations, which leads to the weak search ability of the algorithm in the early stage and the weak convergence in the later stage? For this reason, the study proposes a nonlinear control factor adjustment strategy, as shown in Eq. (13).

$$a = 2\sqrt{1 - \left( t/t_{\max} \right)^2} \quad (13)$$

In Eq. (13), the maximum number of iterations is denoted by $t_{\max}$. The research proposed strategy and the traditional strategy are shown in Fig. 5.

It can be seen that under the strategy proposed in the study, the $a$ value changes slowly at the beginning of the iteration and the search performance of GWO is strong; at the later part of the iteration the $a$ value changes faster, which makes GWO have good convergence. As the GWO algorithm has the disadvantage of maturing and convergent too early, resulting in poor search accuracy, the study proposes a segmentation step strategy that adjusts the update mechanism according to the $A$ value. When $|A| \leq 1$ is used, the update mechanism shown in Eq. (7) is adopted. When $|A| > 1$, three individuals are randomly selected $r_1, r_2, r_3$ to determine the search range of search individuals $R'$. With this tactic, it is possible to provide some of the members of the badly located grey wolves an opportunity to take part in the algorithm's location update choice, successfully increasing population diversity and preventing the early GWO algorithm occurrence. The above can be represented as Fig. 6.

Fig. 5.    Control factor adjustment strategy.



Fig. 6.    Segmental update strategy of GWO algorithm.

Comprehensive above, build MGWO, use MGWO to find the best parameters of FPID, optimize it, build MGWO-FPID model, realize the high precision intelligent temperature control of the instrument, improve the temperature control effect, and then improve the industrial production efficiency, which has positive significance to the safety guarantee of industrial production.

## IV.    PERFORMANCE ANALYSIS OF MGWO-FPID INSTRUMENTATION INTELLIGENT TEMPERATURE CONTROL MODEL

To validate the capability of the MGWO-FPID instrumentation intelligent temperature control model proposed in this study, the study selects an instrumentation used for thermal power metering for the experiment and uses the software Matlab for the simulation experiment. The experimental parameters are as follows, the voltage is 10V, the network signal state temperature, the program state is offline, the operation model is selected as the IoT operation model, and the analysis method uses platform computing analysis. The current state-of-the-art PID control algorithms are particle swarm optimization FPID (PSO-FPID) model and FPID (ISOA-FPID) model optimized based on improved seeker optimization algorithm (ISOA). It collected a large amount of historical temperature data, including temperature changes under different working conditions. Use platform computing

analysis methods to preprocess, clean, and analyze data to extract key features and trends. In the comparison model, Particle Swarm Optimization FPID (PSO-FPID) model: based on particle swarm optimization algorithm, PID parameters are optimized to improve the response speed and stability of temperature control. FPID (ISOA-FPID) optimization model based on improved seeker optimization algorithm (ISOA): By improving the seeker optimization algorithm to adjust PID parameters, the robustness and adaptability of the control are improved. Under the same experimental conditions, run the MGWO-FPID model, PSO-FPID model, and ISOA-FPID model separately. Record the temperature control effects of each model under different working conditions, including control accuracy, response speed, and stability indicators. By comparing the experimental results, analyze the superiority and performance characteristics of the MGWO-FPID model compared to other models. In the experiment, MGWO-FPID was compared with PSO-FPID and ISOA-FPID models to conduct a comprehensive analysis in terms of objective function and fitness, temperature control regulation efficiency, control accuracy, model MAE, recall, and AUC. The temperature control performance of the above three intelligent temperature control models are compared respectively. The particle number of the PSO-FPID model is 45, the inertia weight is 0.8, and the acceleration constant is 1.51; The initial population of the ISOA-FPID model is 45, with a crossover probability of 0.5, a mutation probability of 0.51, and a learning factor of 0.3; The wolf pack size of the MGWO-FPID model, with an initial population of 45, a wolf pack level of 0.55, and contraction and expansion factors of 0.4 and 0.5, respectively. The maximum number of iterations for all models is 300. Firstly, the optimization effects of the above three models are compared. During the iterative process, the changes of the fitness values and the objective function values of MGWO-FPID model, ISOA-FPID model and PSO-FPID model are shown in Fig. 7. It is evident that the MGWO-FPID model's convergence is superior to that of the ISOA-FPID model and the PSO-FPID model because the objective function value of the MGWO-FPID model declines faster and the fitness value increases faster during the iterative process. Compared to the ISOA-FPID model and the PSO-FPID model, the MGWO-FPID model achieves an objective function value of 5 10-8 in Fig. 7(a). This value is four orders of magnitude lower. In Fig. 7(b), the fitness value of the MGWO-FPID model reaches 13.1, which is 2.8 and 3.3 higher than the ISOA-FPID model and the PSO-FPID model, respectively.

The data related to the instrumentation used for thermal power metering were input into the MGWO-FPID model, ISOA-FPID model and PSO-FPID model and simulated using Matlab software, and the simulation curves of several temperature control models are shown in Fig. 8. In Fig. 8, it can be seen that the MGWO-FPID model has no overshoot and completes the temperature control regulation of the instrument in a shorter time compared with the ISOA-FPID model and PSO-FPID model. The above results demonstrate that the MGWO-FPID model is more efficient in temperature control regulation.

(a) Objective function value



(b) Fitness value

Fig. 7.   Changes in fitness values and objective function values of the model.



Fig. 8.   Simulation curves of several temperature control models.

The control accuracies of MGWO-FPID model, ISOA-FPID model and PSO-FPID model in the intelligent control of instrument temperature are shown in Fig. 9. In Fig. 9, it can be seen that the MGWO-FPID model has higher accuracy and requires fewer iterations to reach the best accuracy. 68 iterations are required for the MGWO-FPID model to reach the best accuracy, which is 95 and 137, fewer than the ISOA-FPID model and the PSO-FPID model, respectively. The MGWO-FPID model's control accuracy is 99.52%, which is greater than the ISOA-FPID model's and the PSO-FPID model's, respectively, by 0.53% and 0.77%.

Mean Absolute Error (MAE), also known as Mean Absolute Error, is a commonly used goodness of fit evaluation criterion in regression analysis. It is the average absolute value of the difference between the predicted value and the actual value. Fig. 10 illustrates this comparison between the change in F1 value and the change in MAE during the course of an iteration of the MGWO-FPID model, ISOA-FPID model, and PSO-FPID model. Fig. 10 demonstrated that the F1 values of MGWO-FPID model, ISOA-FPID model and PSO-FPID model are rapidly increasing and the MAE values are rapidly decreasing at the beginning of the iteration, and after the F1 and MAE values reach a certain level, the F1 and MAE values

of MGWO-FPID model, ISOA-FPID model and PSO-FPID model no longer change significantly, indicating that the models have converged. It can be seen that the MGWO-FPID model converges faster. The F1 value of the MGWO-FPID model, which is 0.58% and 1.24% higher than the ISOA-FPID model and PSO-FPID model, respectively, reaches 96.14% as shown in Fig. 10(a). Fig. 10(b) illustrates that the MGWO-FPID model's MAE value is 8.53, which is 1.22 and 2.87 less than the MAE values for the ISOA-FPID model and PSO-FPID model, respectively. The above results can indicate that the MGWO-FPID model has better performance.

Utilising the Recall value as shown in Fig. 11, the performance of the MGWO-FPID model, ISOA-FPID model, and PSO-FPID model is assessed. Fig. 11 illustrates how the MGWO-FPID model has a greater Recall value and better convergence. The Recall value of MGWO-FPID model is 95.37%, which is 0.62% and 1.33% higher than the ISOA-FPID model and PSO-FPID model, respectively.



Fig. 9.   Control accuracy of model in instrument temperature intelligent control.

(a) F1



(b) MAE

Fig. 10. Changes in F1 value and MAE during iteration of the model.



Fig. 11. Recall values for three models.



Fig. 12. ROC curve trends and AUC values for three models.

The comprehensive performance of MGWO-FPID model, ISOA-FPID model, and PSO-FPID model was evaluated by ROC curve trend with AUC value, as shown in Fig. 12. AUC (Area Under the Curve) is a commonly used metric to evaluate the performance of classification models, widely used in fields such as machine learning, data mining, and statistics. The range of AUC values is between 0 and 1, with values closer to 1 indicating better model performance, and values closer to 0.5 indicating relatively random model predictions. As can be noticed, the MGWO-FPID model's AUC value is 0.995, which is 0.011 and 0.024 higher than the AUC values for the ISOA-FPID model and the PSO-FPID model, respectively. The aforementioned results show that the MGWO-FPID instrumentation intelligent temperature control model suggested in the study performs more comprehensively than the other two models. In summary, the MGWO-FPID instrumentation intelligent temperature control model proposed in the study has high accuracy and efficiency, and effective comprehensive performance, which can effectively realize the high-precision intelligent temperature control of the instrumentation, enhance the temperature control effect, and then improve the industrial production efficiency, and has positive significance for the safety guarantee of industrial production.

In order to more intuitively demonstrate the performance of the three models, a composite table was formed based on the above experimental results, as shown in Table I.

TABLE I.    SYNTHETIC TABLE OF PERFORMANCE RESULTS FOR THREE MODELS

| / | MGWO-FPID | ISOA-FPID | PSO-FPID |
|---|---|---|---|
| Fitness value | 13.1 | 10.3 | 9.8 |
| Achieving optimal precision iteration times | 68 | 163 | 205 |
| Control accuracy | 99.52% | 98.99% | 98.75% |
| F1 value | 96.14% | 95.56% | 94.90% |
| MAE | 8.53 | 9.75 | 11.4 |
| Recall value | 95.37% | 94.75% | 94.04% |
| AUC | 0.995 | 0.984 | 0.971 |

## V. RESULTS AND DISCUSSION

The temperature control performance of intelligent temperature control models using three PID control algorithms, MGWO-FPID, PSO-FPID, and ISOA-FPID, was compared in the experiment. Firstly, the optimization effects of three models were compared. In terms of objective function and fitness, the fitness value of the MGWO-FPID model reached 13.1, which is 2.8 and 3.3 higher than the ISOA-FPID model and PSO-FPID model, respectively. In terms of temperature control efficiency, compared with the ISOA-FPID model and PSO-FPID model, the MGWO-FPID model has no overshoot and completes the temperature control adjustment of the instrument in a shorter time. In terms of control accuracy, the MGWO-FPID model requires 68 iterations to achieve optimal accuracy, which is 95 and 137 fewer than the ISOA-FPID model and PSO-FPID model, respectively. The control accuracy of the MGWO-FPID model is 99.52%, which is 0.53% and 0.77% higher than the ISOA-FPID model and PSO-FPID model, respectively. In terms of model MAE, the MGWO-FPID model has a MAE value of 8.53, which is 1.22 and 2.87 lower than the ISOA-FPID model and PSO-FPID model, respectively. In terms of recall rate, the MGWO-FPID model has a recall rate of 95.37%, which is 0.62% and 1.33% higher than the ISOA-FPID model and PSO-FPID model, respectively. In terms of AUC, the MGWO-FPID model has an AUC value of 0.995, which is 0.011 and 0.024 higher than the ISOA-FPID model and PSO-FPID model, respectively. In summary, through experimental comparison, we can conclude that the MGWO-FPID model exhibits better performance in intelligent temperature control compared to the PSO-FPID and ISOA-FPID models. It has significant advantages in terms of objective function and fitness, temperature control regulation efficiency, control accuracy, model MAE, recall, and AUC. Therefore, in actual industrial production, using the MGWO-FPID model for intelligent temperature control will help improve production efficiency and product quality.

## VI. CONCLUSION

In industrial production, the temperature control of the instrument is related to the accuracy of the instrument detection data, which affects the safety and stability of industrial production. Therefore, MGWO-FPID instrumentation intelligent temperature control model is proposed for the current instrumentation intelligent temperature control methods with low accuracy and efficiency defects. Based to the experimental findings, the MGWO-FPID model's objective function value was 510-8, which was 4 orders of magnitude less than that of the ISOA-FPID model and 6 orders of magnitude less than that of the PSO-FPID model; the adaptation degree value reaches 13.1, which is 2.8 and 3.3 higher than ISOA-FPID model and PSO-FPID model respectively; the control regulation time is 2.08s, which is higher than ISOA-FPID model and PSO-FPID model. FPID model and PSO-FPID model, respectively; 68 iterations are required to achieve the best accuracy, which is 95 and 137 times less than the ISOA-FPID model and PSO-FPID model, respectively; the F1 value reaches 96.14%, which is 0.58% and 1.24% higher than the ISOA-FPID model and PSO-FPID model, respectively The MAE value was 8.53, which was 1.22 and 2.87 lower than the ISOA-FPID model and PSO-FPID model, respectively; the Recall value was 95.37%, which was 0.62% and 1.33% higher than the ISOA-FPID model and PSO-FPID model, respectively; the AUC value reached 0.995, which was 0.01% higher than the ISOA-FPID model and PSO-FPID model, respectively. The above results can prove that the MGWO-FPID instrumentation intelligent temperature control model proposed in the study has high accuracy and efficiency, and effective comprehensive performance, which can effectively realize the high precision intelligent temperature control of the instrumentation, improve the temperature control effect, and then enhance the industrial production efficiency, which is of positive significance to the safety guarantee of industrial production. In the experiment, due to limitations in data sources, there may indeed be discrepancies between the experimental results and the actual situation. In order to improve the accuracy and reliability of research, it is indeed necessary to broaden the scope of research to eliminate accidental errors. In order to make the research results more representative, it is necessary to obtain data from a wider range of thermoelectric metering devices to cover a wider range of device performance and possible sources of error; Compare with intelligent temperature control models in other fields to understand their respective advantages and limitations, and further improve the performance of instrument intelligent temperature control models based on MGWO FPID; Factors such as the operator's experience, skills, and psychological factors under environmental conditions can be considered to improve the comprehensiveness of the study.

## FUNDING

## REFERENCES

[1] Uygun A D, Unal M, Falakaloglu S, Guven Y. Comparison of the cyclic fatigue resistance of hyflex EDM, vortex blue, protaper gold, and onecurve nickel -Titanium instruments. Nigerian journal of clinical practice, 2020, 23(1): 41-45.

[2] Alcelay I, Peña E, Omar A A. Hot working behaviour and processing maps of duplex cast steel. International Journal of Materials Research, 2021, 112(7): 518-526.

[3] Gheisarnejad M, Khooban M H. An intelligent non-integer PID controller-based deep reinforcement learning: implementation and experimental results. IEEE Transactions on Industrial Electronics, 2020, 68(4): 3609-3618.

[4] Jabeur C B, Seddik H. Optimized neural networks-PID controller with wind rejection strategy for a Quad-Rotor. Journal of Robotics and Control (JRC), 2022, 3(1): 62-72.

[5] Zhao D, Li F, Ma R, Zhao G, Huangfu G. An unknown input nonlinear observer based fractional order PID control of fuel cell air supply system. IEEE Transactions on Industry Applications, 2020, 56(5): 5523-5532.

[6] Phu N D, Hung N N, Ahmadian A. Senu N. A new fuzzy PID control system based on fuzzy PID controller and fuzzy control process. International Journal of Fuzzy Systems, 2020, 22(7): 2163-2187.

[7] Shi J Z. A fractional order general type-2 fuzzy PID controller design algorithm. IEEE Access, 2020, 8: 52151-52172.

[8] Shi L, Hu Y, Su S, Guo S, Xing H, Hou X, Liu Y, Chen Z, Li Z, Xia D. A fuzzy PID algorithm for a novel miniature spherical robots with three-dimensional underwater motion control. Journal of Bionic Engineering, 2020, 17: 959-969.

[9] Ghamari S M, Narm H G, Mollaee H. Fractional-order fuzzy PID controller design on buck converter with antlion optimization algorithm. IET Control Theory & Applications, 2022, 16(3): 340-352.

[10] Shi Q, Lam H K, Xuan C, Chen M. Adaptive neuro-fuzzy PID controller based on twin delayed deep deterministic policy gradient algorithm. Neurocomputing, 2020, 402: 183-194.

[11] Liu Y, Fan K, Ouyang Q. Intelligent traction control method based on model predictive fuzzy PID control and online optimization for permanent magnetic maglev trains. IEEE Access, 2021, 9: 29032-29046.

[12] Kumar Khadanga R, Kumar A, Panda S. Frequency control in hybrid distributed power systems via type-2 fuzzy PID controller. IET Renewable Power Generation, 2021, 15(8): 1706-1723.

[13] Sain D, Mohan B M. Modeling, simulation and experimental realization of a new nonlinear fuzzy PID controller using Center of Gravity defuzzification. ISA transactions, 2021, 110: 319-327.

[14] Zamfirache I A, Precup R E, Roman R C, Petriu E M. Policy iteration reinforcement learning-based control using a grey wolf optimizer algorithm. Information Sciences, 2022, 585: 162-175.

[15] Liu J, Wei X, Huang H. An improved grey wolf optimization algorithm and its application in path planning. IEEE Access, 2021, 9: 121944-121956.

[16] Hao P, Sobhani B. Application of the improved chaotic grey wolf optimization algorithm as a novel and efficient method for parameter estimation of solid oxide fuel cells model. International Journal of Hydrogen Energy, 2021, 46(73): 36454-36465.

[17] Otair M, Ibrahim O T, Abualigah L, Altalhi M, Sumari P. An enhanced grey wolf optimizer based particle swarm optimizer for intrusion detection system in wireless sensor networks. Wireless Networks, 2022, 28(2): 721-744.

[18] Sun X, Jin Z, Cai Y. Yang Z. Chen L. Grey wolf optimization algorithm based state feedback control for a bearingless permanent magnet synchronous machine. IEEE Transactions on Power Electronics, 2020, 35(12): 13631-13640.

[19] Rajamoorthy R, Arunachalam G, Kasinathan P. Panchal H. Kazem H A. Sharma j P. A novel intelligent transport system charging scheduling for electric vehicles using Grey Wolf Optimizer and Sail Fish Optimization algorithms. Energy Sources, Part A: Recovery, Utilization, and Environmental Effects, 2022, 44(2): 3555-3575.

[20] Xu Z, Yang H, Li J. Zhang X. Lu B. Gao S. Comparative study on single and multiple chaotic maps incorporated grey wolf optimization algorithms. IEEE Access, 2021, 9: 77416-77437.

[21] Ahmadi B, Younesi S, Ceylan O. Ozdemir A. An advanced grey wolf optimization algorithm and its application to planning problem in smart grids. Soft Computing, 2022, 26(8): 3789-3808.

# Designing an Adaptive Effective Intrusion Detection System for Smart Home IoT

## A Device-Specific Approach with SDN Deployment

Hassen Sallay

College of Computing, Umm Al-Qura University, Makkah, KSA

*Abstract*—As the ubiquity of IoT devices in smart homes escalates, so does the vulnerability to cyber threats that exploit weaknesses in device security. Timely and accurate detection of attacks is critical to protect smart home networks. Intrusion Detection Systems (IDS) are a cornerstone in any layered security defense strategy. However, building such a system is challenging given smart home devices' resource constraints and behaviors' diversity. This paper presents an adaptative IDS based on a device-specific approach and SDN deployment. We categorize devices based on traffic profiles to enable specialized architectural design and dynamically assign the suitable detection model. We demonstrate the IDS efficiency, effectiveness, and adaptability by thoroughly benchmarking an ensemble of machine learning models, mainly tree ensemble models and extreme learning machine variants, on the up-to-date IoT CICIoT2023 security dataset. Our IDS multi-component device-aware architecture leverages software-defined networking and virtualized network functions for scalable deployment, with an edge computing design to meet strict latency requirements. The results reveal that our adaptive model selection ensures detection accuracy while maintaining low latency, aligning with the critical requirement of real-time accuracy and adaptability to smart home devices' traffic patterns.

*Keywords*—*Smart home; IoT; IDS; taxonomy, architecture; SDN; ELM*

## I. INTRODUCTION

The rapid growth of the smart home market broadly opens the doors to several threats to people's security and privacy. People are often unaware of security vulnerabilities, and manufacturers fail to prioritize security. This combination leads to a growing attack surface for hackers to exploit. Indeed, it is well known that many smart home devices, including IP cameras, smart locks, smart lighting systems, etc., contain vulnerabilities that attackers can exploit to intrude into home networks. Successful intrusions into IoT devices can allow hackers to not only steal sensitive user data but also take control of critical devices. Hence, there is a growing need for intelligent security systems to detect abnormal behaviors and attacks on smart home IoT devices in real time.

Since no one-size-fits-all security solution exists, a defense-in-depth approach and appropriate design and implementation should be context-aware to protect against threats and specific attack vectors. Among the complementary tools in the security layered defense comes network intrusion detection systems (NIDSs). They are security tools that continuously analyze traffic to identify intrusions and attacks. Traditional IDS employ signature-based detection, which matches known attack patterns. More advanced anomaly detection techniques spot statistical deviations from normal traffic to surface previously unseen attacks. However, building accurate intrusion detection models for IoT is challenging due to several factors. IoT devices have much more resource constraints than traditional computing systems and exhibit complex and dynamic behaviors. Moreover, the constantly evolving threats and vulnerabilities must be efficiently well-tracked for an adaptive security defense in the smart home context. Thus, knowing the ground truth for device and traffic features will be useful in tackling intrusion security challenges posed by smart home environments.

Several research efforts have been spent to tackle these challenges. The research in [1] provided a comprehensive review of intrusion detection systems using machine and deep learning in IoT, discussing challenges, solutions, and future directions. They emphasized the need for efficient and accurate detection methods but did not propose a specific architecture or implementation. The study in [2] surveys network intrusion detection for IoT security based on learning techniques, highlighting the importance of efficient learning algorithms for smart home security. It deeply and thoroughly explores recent works focusing on machine learning techniques. However, its scope does not include architectural design issues such as adaptability and real-time requirements. The study in [3] introduced a deep learning application for invasion detection in industrial IoT sensing systems. While this work is relevant for industrial applications, it may not directly translate to smart home environments due to different operational constraints and attack vectors. The research in [4] proposes an integrated multilayered framework for IoT intrusion decisions and instantiates it for the industrial IoT. Although the framework can be instantiated to the smart home context, the paper did not specify the architectural design and deployment. All these works raised the flag that most existing methods overlook key IoT constraints like low latency, dynamic device behaviors, and resource limitations that impact real-world-scale adoption.

We also cite some works that gave us insights to develop our proposed solution. The study in [5] proposed an intrusion detection system using an Online Sequence Extreme Learning Machine in the advanced metering infrastructure of smart grids. Their model focused on sequential data processing, which is pertinent and can be adapted to the continuous monitoring required in smart homes. The research in [6] discussed intrusion detection in fog computing and Mobile

Edge Computing. This work is particularly relevant as it considers the edge computing paradigm, which is increasingly adopted in smart home IoT. However, it does not discuss the use of machine learning in intrusion detection. The study in [7] presented a self-configurable cyber-physical intrusion detection system for smart homes using reinforcement learning. This system's adaptability to changing conditions in smart homes is a significant step towards dynamic and responsive security systems. The research in [8] explored a hybrid approach using an artificial immune system for intrusion detection in smart home networks. This work highlighted the potential of ELM for fast learning and generalization but did not focus on the real-time aspect of intrusion detection. The study in [9] integrates the software-defined networking, machine learning, and manufacturer usage descriptions standard with an intrusion detection and prevention system to assess its influence on network security. While including standards-based ingredients is interesting, their work is limited to manufacturing and does not consider the architectural effectiveness and network traffic characteristics.

For the traffic characterization, [10] recognizes the importance of understanding IoT data characteristics for modeling the data bursts typical of IoT use cases, and it introduces an advanced ON/OFF traffic modeling approach tailored for the varied applications within a smart city context. While the work is pivotal for statistical modeling of the IoT traffic, it does not consider their solutions' architectural design and deployment. The study in [11] provides insights into IoT traffic characteristics in the specific context of smart home and campus environments. They found that IoT devices exhibit periodic behavior with significant idle time. The devices generate a small amount of traffic, and most communicate with a small number of remote servers, often located in the same country as the device. The study also found several security and privacy issues, including devices communicating over unencrypted channels and devices communicating with servers in countries known for privacy concerns. This paper aims to design and build a flexible, scalable IDS that efficiently ensures security defense in real smart home environments without losing generality and adaptability. The contributions are:

- We introduce a device-aware approach that categorizes IoT devices based on their traffic profiles and behaviors, leading to a more tailored and efficient detection process that can adapt to the heterogeneous nature of smart home devices.

- Our solution employs an ensemble of optimized machine learning models, including extreme learning machine variants chosen based on the device category, balancing the tradeoffs between speed, accuracy, and resource usage.

- We provide an experimental evaluation using an up-to-date security dataset, demonstrating the effectiveness of our approach in a realistic smart home context.

- We design a multi-component IDS architecture using network traffic profiles for real-time intrusion detection in smart home IoT environments. Our architecture leverages software-defined networking (SDN) and virtualized network functions (VNFs), allowing for a flexible and scalable deployment that can be adapted for both cloud and edge computing scenarios. We mainly opted for an edge computing-based deployment on an SDN testbed to meet strict latency requirements.

The remainder of this paper is organized as follows: Section II presents our methodology steps. Section III characterizes the smart home devices' traffic and presents a simple traffic-based taxonomy. Section IV shows the architecture design and deployment. Section V presents the benchmarking of the machine learning models. Section VI shows the benchmarking results. Section VII concludes the paper and gives some future works.

## II. METHODOLOGY

We propose a two-stage methodology within four steps to develop our intrusion detection system:

### A. Devices' traffic characterization

*1)* Explore and categorize the commonly used devices in the smart home environment.

*2)* Characterize the traffic devices and build a traffic-centric devices taxonomy.

### B. Architectural Design and Model Selection

*1)* Design and deployment of a smart home IDS-tailored architecture.

*2)* Benchmark the ML models on a recent dataset and select the appropriate model based on the previous steps.

More specifically, we start by enumerating devices and device/data categories to understand the ecosystem, and then we characterize traffic patterns and classify devices into a useful taxonomy. We are leveraging this knowledge to design and deploy suitable IDS architecture. Then, select appropriate machine learning models that detect intrusions optimized for the specifics of the smart home domain. The result is an IDS purpose-built to the unique smart home environment versus more generic systems. The key rationale is that threats exploit specific device vulnerabilities and traffic flows in the smart home, so an IDS must be aware of these devices and patterns to identify attacks. The proposed methodology builds this intrinsic knowledge by examining the ecosystem to customize the IDS. This context-aware solution can better distinguish attacks from normal traffic and has utility detecting intrusions that more generic learning-based systems may overlook in the IoT setting.

## III. DEVICES' TRAFFIC CHARACTERIZATION

A typical smart home would include various IoT devices commonly used, such as smart thermostats, smart lighting systems, smart security cameras, smart locks, smart appliances (e.g., refrigerators, washing machines), smart speakers or home assistants, and smart TVs. Table I characterizes common consumer IoT devices along three dimensions: subcategories based on features, key devices' behavior patterns, and data reflecting network traffic patterns in size and time. Grouping into broader categories like smart speakers while still enumerating specific device types enables roll-up

summarization and device-specific analysis. Smart TVs, for instance, have subcategories of basic models focused on video streaming with basic controls versus smart TVs, which have an integrated app platform enabling third-party applications like Netflix and YouTube. The use cases cover on-demand video and accessing these apps for entertainment and information. This expanded functionality versus standalone streaming leads to more diverse, multimodal traffic encompassing control commands, actual video data, and app communications.

Referring to Table I and the data nature, we see a mix of primarily unimodal control traffic for simpler devices like lightbulbs, thermostats, and locks with relatively straightforward command and monitoring use cases. Meanwhile, more advanced devices like cameras, speakers, and fridges demonstrate bimodal or multimodal traffic indicative of more mixed media, including audio, video, and firmware downloads, resulting in variable network utilization.

The complexity arising from supporting multiple integrated apps paired with streaming video in one device results in a complex traffic profile that may require advanced analytic approaches beyond simple machine learning algorithms to adequately characterize if simple models prove to have insufficient descriptive capability and predictive accuracy. However, for unimodal traffic, simpler models should suffice without overcomplicating analysis.

Thus, we built a simple taxonomy of smart home IoT devices based on their network traffic characteristics:

- Streaming Devices [Smart speaker, IP camera, Voice assistant robot] (Traffic patterns are: (1) Bimodal packet size distribution (small control + large streaming packets), (2) Bursty packet timing during streaming, and (3) Higher and variable traffic volume.)

- Intermittent Control Devices [Smart lightbulb, Smart thermostat, Smart fridge, Smart doorbell, Smart blinds, Irrigation controller] (Traffic patterns are: (1) Uniformly small packet sizes, (2) Periodic keepalives + event-driven commands, and (3) Low traffic volume with occasional spikes )

- Monitoring Devices [Smoke detector, Motion sensor, Door/window sensor] (Traffic patterns are: (1) Small packets for status updates, (2) Sporadic or periodic timing, and (3) Very low traffic volume)

- Actuators [Smart lock, Garage door opener, Smart plug] (Traffic patterns are: (1) Small command packets, (2) Event-driven timing (3) Extremely low traffic volume).

Based on this taxonomy, considering device behaviors and traffic profiles, we categorize home IoT devices into two broad classes for architectural design:

- High-throughput devices include video cameras, media hubs, etc., generating high volumes of multimedia traffic. The patterns are more complex and variable.

- Low-throughput devices consist of simpler sensors and controllers for lighting, smoker detectors, etc., with minimal traffic. The patterns tend to be regular and predictable.

Accordingly, in the next section, we propose a multi-component IDS architecture to secure the smart home.

TABLE I.     DEVICES TRAFFIC CHARACTERISTICS

| Device Category | Device | | Data | |
|---|---|---|---|---|
| | *Subcategory* | *Behavior* | *Size* | *Timing* |
| Smart TV | Basic, Smart | Intermittent streaming | Bimodal (control + audio packets) | Bursty during use and periodic otherwise |
| Smart Speaker | Audio streaming, Voice assistant, Smart display | Mostly control commands | Uniformly small | Event-driven/periodic |
| Smart lightbulb | Tunable white, RGB, Motion sensor | Continuous video | Bimodal (small + large video packets) | Periodic real-time streaming |
| IP camera | Video doorbell, Baby monitor, Security camera | Infrequent controls | Uniformly small | Periodic polling + event-driven |
| Smart thermostat | Self-contained, HVAC integrated | Intermittent traffic | Bimodal (control + firmware updates) | Periodic sensors updates |
| Smart fridge | Display model, Bottom-freezer model | Sparse controls | Uniformly small | Infrequent periodic keepalives |
| Smart lock | Bluetooth, WiFi, Z-Wave | Activated when used | Small control packets | Event-driven only |
| Garage door opener | WiFi/Bluetooth connected, Remote controlled | Intermittent controls | Small control packets | Periodic + event-driven |
| Smart blinds | Motorized, App/voice controlled | Sparse status report | Small power toggling packets | Periodic status updates |
| Smart plug | Controllable, Monitored | Event-driven alerts | Small alert packets | Sporadic alarms, periodic heartbeats |
| Smoke detector | Integrated, Smart alarm | Regularly scheduled operation | Small control/status packets | Periodic polling + daily schedules |
| Irrigation controller | App connected, Weather adjusted | Intermittent streaming | Packet Size Distribution | Packet Timing Distribution |

## IV. ARCHITECTURAL DESIGN AND MODEL SELECTION

### A. Architecture Design

The core of our proposed intrusion detection system comprises a multi-component architecture tailored to secure diverse IoT devices in smart home environments. Our design meets the following key requirements: real-time detection capability, adaptability to evolving behaviors, detection accuracy for known and zero-day attacks, and computational efficiency to operate given smart home resource constraints.

The modular architecture allows customizing specific components to address deployment-specific needs. The IoT devices would commonly be connected to a local network, including a gateway or router, to manage network traffic and connect to the Internet. The network commonly includes a mix of wired and wireless connections, depending on the specific devices used, and HTTP(S), MQTT, CoAP, and Zigbee are the common IoT-used protocols. Furthermore, smart home often has a network firewall and other basic security measures for each device the manufacturer provides. The security threats landscape includes denial or distributed denial of service (DoS/DDoS) attacks, malware or ransomware attacks, unauthorized access or intrusion attempts, and data breaches or exfiltration attempts. We add a layer for intrusion detection that stays behind the firewall. Mainly the NIDS system includes the following components:

- Traffic Inspector capturing and pre-processing all device traffic flows. It mainly (1) captures raw network traffic using port mirroring, (2) extracts flow-based features like source/destination IP, ports, packet sizes, etc., (3) tags flow with device identities from logs, and (4) forwards processed flows to Device Profiler.

- Device Profiler identifies and assigns device type to a high/low throughput category. It mainly (1) maintains an inventory of identified IoT devices, (2) classifies devices into high or low throughput groups, and (3) pushes device type and group to ML Model Selector.

- ML Model Selector chooses the optimal intrusion detection model for that device type. It mainly (1) houses a catalog of optimized ML models for each device group, (2) models tailored for the complexity and behaviors of that group, (3) queries device group for a flow from Device Profiler and (4) dynamically it selects the matching model for anomaly detection.

- Model Repository contains specialized ML models tailored for each device class. It mainly (1) stores specialized ML models, (2) contains different algorithms that suit traffic complexities, and (3) includes models pre-trained on normal and attack device data.

- Intrusion Detector to analyze traffic for intrusion using a selected model. It mainly (1) receives network traffic flow features, (2) feeds to Model Selector chosen model, (3) the model analyzes the sequence for intrusions, and (4) the classifier flags intrusion if found.

- Alert Manager raising intrusion alerts as needed with attack details. It mainly (1) collects intrusion alerts from Intrusion Detector, (2) provides details like affected device attack type, and (3) raises notifications to admin and response systems.

We first utilize a Traffic Inspector module that captures raw network traffic using port mirroring techniques. It then extracts flow-based features like source and destination IPs, ports, packet sizes, and tag flows to specific device identities obtained from logs. The processed traffic flows are forwarded to a Device Profiler component, which maintains an inventory of devices identified on the network. Leveraging both domain knowledge, the Device Profiler categorizes devices into either high throughput or low throughput groups. High throughput devices like cameras and media hubs generate higher volumes of multimedia network traffic with more complex and variable patterns. In contrast, simpler sensors and controllers constitute the low throughput group with minimal and regular traffic.

The device type and group information are passed into an ML Model Selector module that maintains a catalog of specialized models tailored for each device group. When the Model Selector receives a query with the device group for a particular traffic flow, it dynamically selects the matching specialized model to analyze that flow for intrusions. This model repository containing diverse algorithms suited for varying traffic complexities is pre-trained on normal and attack data generated from devices in the corresponding category.

An Intrusion Detector module takes the network traffic flow features and feeds them into the model instance chosen by the Model Selector for that flow. Based on previous learning, the selected model analyzes the sequence to detect intrusions, finally flagging likely security intrusions. Any intrusion alerts are collected by an Alert Manager, who provides details like the affected device and attack type to administrators and incident response systems.

### B. Architecture Deployment

Our proposed intrusion detection system's components leverage software-defined networking (SDN) capabilities for efficient and flexible system deployment [12,13]. The SDN controller provides a central orchestration point for the various IDS modules [14]. Network switches are configured using SDN policies to mirror IoT traffic flows that need to be inspected, tapping them to feed into the IDS Traffic Inspector module. The centralized network view within the SDN control plane also enables mapping these flows to specific IoT devices on the network.

A software-defined implementation offers significant advantages in flexibility, programmability, and scalability. The centralized control plane greatly simplifies tapping into a high volume of IoT flows in dynamic environments while automating complex policy configurations needed for mirroring. Device profiles and policies can be updated easily as new IoT devices get added over time. SDN also enables large-scale deployments with intelligent traffic engineering and usage optimization across available IDS resources. Therefore, an SDN-based deployment for the intrusion detection

infrastructure makes our IDS more agile and adaptive, mainly as smart home IoT adoption grows exponentially.

Practically, the core detection modules of the IDS, including the Device Profiler, ML model Selector, and Intrusion Detector, are implemented as virtualized network functions (VNFs). By leveraging VNFs and placing them flexibly on commodity servers, we scale out these modules on demand to meet the throughput needs of real-time detection across many IoT devices. As device diversity expands or new models are added to the model repository, more VNF instances can be spun up accordingly. The global view allows the SDN controller to intelligently load balance traffic flows across the VNF resources for optimal efficiency.

The VNF-based deployment can leverage both cloud and edge computing approaches: (1) Cloud-based Deployment where the VNFs for the IDS components like Traffic Inspector, Device Profiler, Model Repository, and Intrusion Detector can be hosted on virtual machines or containers in a private or public cloud. This allows leveraging cloud platforms' flexibility, scalability and managed services. The globally distributed nature of major cloud providers also allows VNFs to be placed closer to IoT deployments for lower latency. However, wide-area network traffic and cloud usage costs may be concerns. (2) Edge Computing Deployment, where we deploy the VNFs on edge servers directly located in smart homes. Edge computing overcomes cloud-based analysis's latency and bandwidth challenges by processing data locally. It provides better responsiveness for real-time intrusion detection [15, 16]. Edge servers can also interface with hardware accelerators for efficient ML model inference. While cloud and edge are viable deployment options, edge computing is better aligned to meet the low latency requirements for real-time intrusion detection across smart home installations. Indeed, the proximity of edge servers to IoT environments makes the IDS more adaptive.

The deployment experimentation can be performed by an SDN testbed where we integrate the edge computing-based deployment with the Mininet/Ryu [17]. Mainly, we set up edge computing nodes in the Mininet topology to host the VNFs (Device Profiler, Model repository, Intrusion Detector). These would consist of lightweight Docker containers. Then, we configure the Ryu controller to steer copies of IoT traffic flows to the nearest edge node for intrusion detection analysis. This mimics real-world edge deployment. The VNF containers process the mirrored device traffic, generate alerts if needed, and export IDS telemetry data. We expose the VNFs via REST APIs for integration with the Ryu controller and monitoring software and evaluate overall latency from the IoT devices to the edge-based IDS VNFs during attack scenarios in Mininet. We can then analyze the responsiveness, overhead, and accuracy relative to an Edge-based deployment. This deployment allows prototyping and demonstrating the benefits of edge computing for IoT environments, leveraging Ryu's programmability and Mininet's flexibility. Automated traffic steering to nearby edge nodes also validates the low latency premise.

## C. Architecture Suitability

Following, we discuss how the proposed modular multi-components IDS architecture design and its SDN-based deployment, along with the Edge-computing technology, help to meet key requirements of real-time detection, adaptability, and accuracy:

*1) Real-time detection capability*: The lean and specialized machine learning models ensure low latency between packet capture by the Traffic Inspector and intrusion alert generation by the Intrusion Detector. In the next section, we will show that the selected models are optimized for efficiency without sacrificing detection accuracy. The virtualized deployment also allows dynamic scaling of detection modules to match incoming traffic volumes. Together, these allow the IDS to provide real-time, sub-second analysis of IoT traffic flows to meet real-time detection needs.

*2) Adaptability to evolving behaviors*: The feedback loop from the Device Profiler to the ML Model Selector allows the system to adapt to changes in device behaviors over time. As traffic patterns change, updated device profiles trigger selection of different models tailored to new behaviors. The models themselves, through re-training, will also adapt during operational use as they observe more data. This tight integration between device knowledge and flexible model selection allows the IDS to adjust to evolving IoT environments.

*3) Detection accuracy*: The model repository for the device category allows highly accurate intrusion detection based on specific device profiles. Tailoring models to capture different IoT devices' normal/attack behavior patterns results in a solution that outperforms one-size-fits-all approaches.

## V. INTRUSION DETECTION MODELS BENCHMARKING

### A. Methodology and Dataset

As per our proposed methodology, we leverage the CICIoT2023 dataset in [18] to categorize smart home IoT devices based on network traffic profiles and select suitable ML models for intrusion detection accordingly. The CICIoT2023 dataset has been created to accelerate research into security analytics and intrusion detection systems tailored for smart home IoT environments. It contains network traffic captures from an extensive smart home IoT testbed comprising over 100 heterogeneous devices. The key value of CICIoT2023 lies in the 33 contemporary IoT-focused attacks spanning seven categories executed on the devices. Many of these attack types are unavailable in other IoT IDS datasets. The attacks leverage compromised IoT devices to penetrate the network, enabling the evaluation of multi-vector IoT threats. Therefore, CICIoT2023 is an invaluable, up-to-date resource for further research into robust intrusion detection tailored for smart home IoT environments facing escalating threats. Moreover, the CICIoT2023 experiments provide data-driven guidance for model selection in our multi-component IDS architecture that analyzes runtime traffic profiles to pick the optimal intrusion detection model tailored to IoT device behaviors and capacities.

The CICIoT2023 is an unbalanced dataset since it is attack intensive. We deploy the SMOT technique to generate a balanced dataset sample. Then, we subsample each dataset into low and high-traffic flow groups that align with our taxonomy of simple and complex IoT devices, respectively. The CICIoT2023 dataset contains 46 features describing the traffic flows. The flow rate in packets/second feature is the most discriminative feature for this clustering across all the available dataset attributes. Then, we experiment with various models on both traffic datasets (balanced and unbalanced), measuring evaluation and model efficiency metrics. Comparing these metrics reveals how different algorithms fare in detecting IoT intrusions for simple and complex device categories, respectively. Additionally, we benchmark the models on the full dataset to validate the improvements gained from our approach of tailored model selection per device traffic profiles. By correlating the model evaluation and efficiency metrics, we can determine the detection accuracy and precision vs. latency tradeoff and evaluate our meeting degree to our key IDS objectives around real-time alerts, adaptability to varying traffic volumes and attack types, and detection precision for IoT environments within typical resource constraints.

*B. Models Selection*

*1) Intuitive analysis of models' suitability:* We intuitively discuss here the suitability of the ML models for the IDS requirements in our context. Tree-based models construct multiple shallow decision trees. They capture nonlinear interactions and complex patterns like network attacks through branching decisions; the tree ensembles balance bias and variance. The tree architectures also suit evolving data through continuous model updates and handheld high-dimensional network data. Support vector machine (SVM) is known for its effectiveness with high dimensional multimodal data, but its complexity impacts real-time performance. While Deep learning models, based on many neural networks, identify complex patterns and capture sequential dependencies, helping detect multi-stage attacks, they require significant data and computing resources that may not suit resource-constrained IoT context. One interesting approach to reducing this complexity is the Extreme Learning Machine (ELM) [11], a fast, single-layer feedforward neural network for classification and regression. They randomly initialize input layer weights and analytically determine output weights. This allows very fast model training suited for real-time usage.

Intuitively, the tree ensemble and ELM class of models seem optimal for balancing efficiency, accuracy, and adaptability within typical IoT constraints. The modular architecture enables deploying complex models like support vector machine or deep learning techniques selectively for capable devices while using tree ensemble and ELM for most real-time detection. The Model Selector dynamically handles this model assignment per device profile. Based on this intuitive analysis, we selected the following set of ML models for benchmarking. Following is a brief description of each model:

*2) Benchmarked models:* We select 12 models to perform our benchmark. Six of them are variants of the ELM, which intuitively brings a promising adequacy for our design requirements. We experiment with its variants [19-23]:

- Kernel ELM is an extension of basic ELM that applies kernel functions like sigmoid, radial basis function (RBF), hyperbolic tangent (tanh), etc., to non-linearly map the input data to new feature spaces before output weight computation. This adds nonlinearity to improve model learning capability for complex patterns.

- Regularized ELM imposes additional constraints on optimizing the output weights matrix calculation. Regularization parameters control model complexity to prevent overfitting, enabling more robust intrusion detection.

- Weighted ELM introduces random scaling factors or weights when multiplying the input layer feature values during forward propagation. This acts as a regularizer like dropout techniques in neural networks, reducing inter-dependencies and improving generalizability.

- OS-ELM (Online Sequential ELM) is the online sequential version of ELM that processes streaming data instance-by-instance for model updates rather than batch learning. This fast incremental learning allows continuous real-time model adaptation, useful for evolving traffic in our context.

- Voting ELM is an ensemble method that trains multiple OS-ELM models on bootstrap samples of the original data. During prediction, the OS-ELM outputs are aggregated through voting to output the overall class.

- Bagging ELM is another ensemble technique using bootstrap sampling to train multiple OS-ELM models. The predictions are aggregated by weighted averaging rather than voting.

- The other six used ML models in the benchmarking are [24-27]:

- Logistic regression is a linear classification model that assigns probabilities to data points belonging to classes using the logistic/sigmoid function. It is fast to train but assumes linear decision boundaries.

- Decision Tree is a simple hierarchical model with branching decisions based on feature thresholds. It is interpretable but prone to overfitting with noisy IoT data.

- Random Forest is an ensemble method combining predictions from many uncorrelated decision trees. It averages out bias and variance for robust performance despite some complexity.

- AdaBoost is another ensemble technique that iteratively focuses on misclassified instances. It can reduce bias and variance errors despite the complexity cost.

- XGBoost is a scalable tree-boosting system for both classification and regression. It uses regularized model formalization for controllable complexity and prevents overfitting.

- LightGBM is a gradient-boosting framework specifically engineered for efficiency, performance, and lower memory usage.

### C. Performance Metrics

- We choose the following key evaluation metrics for assessing our benchmark models' performance:

- Training Time (s): selecting efficient models is crucial for real-time model updates as new devices get added and data change in volume and nature. Therefore, lower training time is better.

- Prediction Time (s) measures classifying flow instances. It impacts the real-time intrusion alert generation capability. Therefore, a lower prediction time is better.

- Latency Time (s) is the sum of training and prediction times reflecting the end-to-end delay from data ingestion to producing an alert. This metric is directly relevant for real-time adaptive detection needs.

- Memory Usage (MB) measures the RAM consumed during model training and inference. It is important due to memory constraints in embedded IoT devices.

- Accuracy is the fraction of correctly classified instances. It provides an overall performance measure but can be misleading for imbalanced data.

- Precision is the fraction calculated by dividing True Positives over the sum of True Positives and False Positives. It is important to minimize false alarms which disrupt users.

- Recall is the fraction calculated by dividing True Positives by the sum of True Positives and False Negatives. It is crucial to maximize the detection of actual attacks and intrusions.

- F1 score is a harmonic mean of precision and recall balancing both metrics. It provides a good measure of detection capability.

Analyzing the latency time tradeoff with accuracy and other metrics allows us to determine the models most suited to real-time, adaptive intrusion detection requirements in smart home IoT environments.

## VI. RESULTS AND DISCUSSION

The results presented for the various machine learning models indicate a diverse range of performance outcomes, with particular interest in the balance between latency and detection metrics such as accuracy and precision. Since we have two datasets (unbalanced and balanced), the figures show the precision metric related to latency times for the unbalanced dataset since the accuracy can be misleading for unbalanced datasets. This is in contrast to the balanced dataset, where accuracy is appropriate. We generated 8*3*2 plots for the 12 models since we measured 8 metrics (training time, prediction time, latency time, memory usage, accuracy, precision, recall, and F1 score) for 3 datasets (low rate, high rate, and low + high rate) sampled from 2 datasets (unbalanced and balanced). Hereafter, we show a subset of the plots most related to the IDS architecture design requirement: real-time, accuracy/precision, and adaptability.

### A. Devices with Low-rate Traffic

For the unbalanced CICIOT2023 dataset, Fig. 1 and Fig. 2 show that the Extreme Learning Machine algorithms demonstrate fast training and prediction latency times, meeting real-time intrusion detection needs. Specifically, the Regularized ELM provides the best balance with strong precision and total latency time. Decision Tree has low latency, making it a good candidate for low-rate devices, as intuitively predicted. In contrast, ensemble methods can enhance precision but have high training times.

For the balanced dataset, the Regularized ELM has the shortest training time at 1.0281 seconds, contributing to a low overall latency time of 1.3915 seconds when combined with its prediction time. This suggests that Regularized ELM is highly suitable for real-time applications where quick model training and prediction are critical. However, the regulated ELM has high precision and recall while maintaining 99.7% accuracy, slightly lower than other models. This combination of speed and reliability makes it the top choice if minimizing detection delay is critical. On the other end of the spectrum, ensemble methods like Random Forest, AdaBoost, XGBoost, and LightGBM have significantly longer training times, contributing to their higher overall latency times making them less suitable for real-time intrusion detection.

Furthermore, the memory usage across all ELM models is relatively smaller than the other models. An interesting outlier is the Decision Tree model, which has low latency similar to ELM methods. The models are fairly consistent regarding memory usage, ranging from 1-1.1 GB, which is acceptable for Edge-computing intrusion detection architecture.

### B. Devices with High-rate Traffic

Fig. 3 and Fig. 4 show that the ELM variants (OS-ELM, Kernel ELM, Regularized ELM, Weighted ELM) have very low latency times, under two seconds. This makes them well-suited for real-time intrusion detection, where getting alerts quickly is critical. However, their detection accuracy is slightly lower than that of ensemble methods, which are slightly accurate but come at the cost of higher latency times. This suggests that Regularized ELM is particularly well-suited for environments where rapid detection is critical and resources may be limited, which fits our case. In contrast, tree ensemble methods show higher latency times, with AdaBoost reaching up to 15.9998 seconds. However, these models exhibit excellent detection performance. The tradeoff is between the higher computational and time costs against the benefit of potentially more accurate detection.

Fig. 1.    Low-rate device training, prediction, latency times, and memory usage for unbalanced and balanced datasets.



Fig. 2.    Latency time according to precision (resp. accuracy) for unbalanced (resp balanced) datasets.

Fig. 3.    High-rate device training, prediction, latency times, and memory usage for unbalanced and balanced datasets.



Fig. 4.    Latency time according to precision (resp. accuracy) for unbalanced (resp balanced) datasets.

The ELM variants have the lowest training and prediction times for the balanced dataset. This enables very fast intrusion detection, suitable for real-time intrusion alerting. Regularized ELM stands out with high accuracy, precision, recall, and F1 scores at 99% despite having the lowest latency of just 1.44 seconds. This demonstrates it can deliver both speed and accuracy. The ensemble methods achieve nearly perfect evaluation metrics but with prohibitive training times. Their high latency makes them impractical for time-critical detection. Memory usage hovers between 1.1-1.3GB for most models with the lowest value of ELM variants. This confirms the adequacy of complex models for the high-rate devices reflecting the complexity of the traffic pattern.

*C. Devices with Low+high Rate Traffic*

Fig. 5 and Fig. 6 show that the Regularized ELM offers an impressive balance between speed and performance, with the lowest training time of 0.9867 seconds and a very competitive prediction time of 0.3748 seconds, resulting in a total latency of 1.3615 seconds. This is complemented by high accuracy

(0.9986), precision (0.9989), and an F1 score (0.9993), making it a strong candidate for real-time applications where both speed and accuracy are critical. In contrast, ensemble methods exhibit higher latency times, ranging from 7.9752 to 31.2592 seconds, but with high accuracy and precision.

The memory usage metric indicates that all models are relatively memory-intensive, with usage ranging from 1601.1875 to 1720.7383 MB. This is a limiting factor in resource-constrained environments, such as our edge computing, and confirms our architecture's suitability, which separates the traffic into two groups to master the time and memory usage consumption.

The training times align with overall latency, with the ELM models being the fastest to train while the ensemble methods are slower. Prediction times are consistently low across the models. Memory usage hovers between 1.6-1.7GB for most models, up to 1.72 GB for LightGBM. So, memory is unlikely to be a constraint.



Fig. 5. High-rate device training, prediction, latency times, and memory usage for unbalanced and balanced datasets.

Fig. 6. Latency time according to precision (resp. accuracy) for unbalanced (resp balanced) datasets.

For the balanced dataset, the ELM variants (OS-ELM, Kernel ELM, Regularized ELM) deliver the fastest performance with training and prediction times under 1 second, maintaining near-perfect accuracy. The low latency makes them optimal for real-time usage. Conversely, ensemble methods like Random Forest, AdaBoost, XGBoost, and LightGBM achieve near-perfect scores on all metrics but at the cost of very high training times, from 137 to over 248 seconds. Their latency is prohibitive for time-critical detection.

Memory usage is reasonably consistent for most models between 1.8 GB and 2GB, which could be limiting for memory-constrained environments. The extreme learner methods have lower memory requirements that may suit resource-limited IoT devices better.

## VII. CONCLUSION AND FUTURE WORK

In this paper, we built a simple taxonomy to separate devices into two broad categories: one with high-volume and complex patterns and low-rate and simple traffic patterns. Based on this categorization, we design a suitable NIDS architecture and select machine learning models that are most adequate for each device type and traffic profile. The models are chosen based on detection accuracy, computational efficiency, and ability to handle complex traffic patterns. For that, we leverage the new CICIoT2023 dataset containing up to date IoT network traffic data with different realistic attacks. Using this dataset, we evaluate various machine learning models to develop an IDS focused on real-time, adaptive detection of intrusions specific to IoT devices.

We optimized the intrusion detection capabilities across smart home IoT deployments by matching the right models to the specific use case requirements around timing, accuracy, and resource constraints. The benchmark analysis guides on selecting between fast ELM variants versus slower but more precise ensemble methods. For low throughput IoT devices with minimal, regular traffic patterns, simpler models like

decision trees provide efficient and fast anomaly detection. Their basic architectures allow quick training and scoring to enable real-time intrusion alerting. In contrast, high throughput multimedia devices require more advanced models like Regularized ELM to capture complex and evolving traffic patterns while maintaining low latency. The nonlinear mappings and optimized complexity in Regularized ELM balance speed and accuracy. So, the device traffic profiles and characteristics directly inform the machine learning model selection to optimize detection capabilities. The benchmark analysis maps models to the specific performance needs driven by the IoT taxonomy of low and high throughput groups. This specialized, profile-based model assignment enhances both efficiency and security.

As a future work, we intend to deploy and evaluate the system in a real-world smart home environment at scale to assess performance with live traffic and attacks. We will also enhance the capability of the system to become intrusion detection and prevention system by correlating intrusion alerts with device vulnerabilities and risk profiles to harden IoT device configurations through SDN dynamically.

## REFERENCES

[1] J., Asharf, N., et all. A review of intrusion detection systems using machine and deep learning in Internet of things: Challenges, solutions and future directions. Electronics, vol. 9 no. 7, pp. 11-77, 2020.

[2] N., Chaabouni, M., Mosbah, A., Zemmari, C., Sauvignac, and P. Faruki, Network intrusion detection for IoT security based on learning techniques. IEEE Communications Surveys & Tutorials, vol. 21. no. 3, pp. 2671-2701. 2019.

[3] G. Altan, SecureDeepNet-IoT: A deep learning application for invasion detection in industrial Internet of things sensing systems. Transactions

on Emerging Telecommunications Technologies, vol.32, no 4, pp. 42-28.2021.

[4]    H. Sallay, An integrated multilayered framework for IoT security intrusion decisions. Intelligent Automation & Soft Computing, vol 36. no 1.pp. 429-444.2023.

[5]    Li, Y., Qiu, R., & Jing, S. Intrusion detection system using Online Sequence Extreme Learning Machine (OS-ELM) in advanced metering infrastructure of smart grid. PloS one, vol. 13. no 2. 2018.

[6]    X., An, X., Zhou, X., Lü, F., Lin, and L. Yang. Sample selected extreme learning machine based intrusion detection in fog computing and MEC. Wireless Communications and Mobile Computing, pp.1-10.2018.

[7]    R., Heartfield, G., Loukas, A., Bezemskij, and E. Panaousis, Self-configurable cyber-physical intrusion detection for smart homes using reinforcement learning. IEEE Transactions on Information Forensics and Security, vol. 16, pp. 1720-1735. 2020.

[8]    E. D. Alalade. Intrusion detection system in smart home network using artificial immune system and extreme learning machine hybrid approach. In 2020 IEEE 6th World Forum on Internet of Things (WF-IoT). pp. 1-2. 2020.

[9]    M., Noman S., Rosli A.H., Mohammad and, Z. Muhammad. SDN based intrusion detection and prevention systems using manufacturer usage description: a survey. (IJACSA) International Journal of Advanced Computer Science and Applications, vol 11, no 12, 2020.

[10]   A. S., Ibrahim, K. Y., Youssef, H., Kamel, and M. Abouelatta. Traffic modelling of smart city internet of things architecture. IET Communications, vol 4, no 8, pp.1275-1284.2020.

[11]   S., Kumar, et all. Characterizing IoT traffic in smart home and campus environments. Proceedings IEEE INFOCOM pp.2706-2715.2020.

[12]   N., Feamster, J., Rexford, and E. Zegura, The road to SDN: an intellectual history of programmable networks. ACM SIGCOMM Computer Communication Review, vol. 44, no. 2, pp. 87-98, 2014.

[13]   K., Benzekki, A., El Fergougui, and A. Elbelrhiti Elalaoui. Software defined networking (SDN): a survey. Security and communication networks, vol. 9, no. 18, pp. 5803-5833.2016.

[14]   L., Mamushiane, A., Lysko, and S. Dlamini,. A comparative evaluation of the performance of popular SDN controllers. In 2018 Wireless Days (WD) (pp. 54-59). 2018.

[15]   F., Liu, et all. A survey on edge computing systems and tools. Proceedings of the IEEE, vol. 107, no. 8, pp.1537-1562.2019.

[16]   S., Hamdan, M., Ayyash, and S. Almajali, Edge-computing architectures for Internet of things applications: A survey. Sensors, vol. 20, no. 22, 6441.2020.

[17]   D., Dholakiya, T., Kshirsagar, and A. Nayak, Survey of mininet challenges, opportunities, and application in software-defined network (sdn). Information and Communication Technology for Intelligent Systems: Proceedings of ICTIS 2020, vol 2, pp. 213-221.2021.

[18]   E.,Neto, et all. CICIoT2023: A real-time dataset and benchmark for large-scale attacks in IoT environment. Sensors, vol. 23, no. 13, 5941, 2023.

[19]   G. B., Huang, Q. Y., Zhu, and C. K. Siew, Extreme learning machine: theory and applications. Neurocomputing, vol. 70, no 1-3, pp. 489-501.2006.

[20]   G. B., Huang, D. H., Wang, and Y. Lan, Extreme learning machines: a survey. International journal of machine learning and cybernetics, vol. 2, pp. 107-122. 2011.

[21]   W., Deng, Q., Zheng, and L.Chen, Regularized extreme learning machine. In 2009 IEEE symposium on computational intelligence and data mining. pp. 389-395. 2009.

[22]   G. B., Huang, S., Song, and K. You, Trends in extreme learning machines: A review. Neural Networks, vol. 61, pp. 32-48.2015.

[23]   O. A., Alade, A., Selamat, and R.Sallehuddin. A review of advances in extreme learning machine techniques and its applications. In Recent Trends in Information and Communication Technology: Proceedings of the 2nd International Conference of Reliable Information and Communication Technology. pp. 885-895. 2018.

[24]   I. D., Mienye, Y.Sun. A survey of ensemble learning: Concepts, algorithms, applications, and prospects. IEEE Access, vol. 10, pp. 129-149. 2022.

[25]   D. Kumar, Priyanka. Decision tree classifier: a detailed survey. International Journal of Information and Decision Sciences, vol. 12, no. 3, pp. 246-269.2020.

[26]   P. A. A., Resende, A. C. Drummond . A survey of random forest based methods for intrusion detection systems. ACM Computing Surveys (CSUR), vol. 51, no. 3, pp. 1-36.2018.

[27]   G., Ke, et all. Lightgbm: A highly efficient gradient boosting decision tree. Advances in neural information processing systems, vol. 30.2017.

# Audio Style Conversion Based on AutoML and Big Data Analysis

Dan Chi

School of Liberal Arts Education and Art Media, Xiamen Institute of Technology, Xiamen, 361000, China

*Abstract*—In the field of audio style conversion research, the application of AutoML and big data analysis has shown great potential. The study used AutoML and big data analysis methods to conduct deep learning on audio styles, especially in style transitions between flutes and violins. The results show that using iterative learning for audio style conversion training, the training curve tends to stabilize after 100 iterations, while the validation curve reaches stability after 175 iterations. In terms of efficiency analysis, the efficiency of the yellow curve and the green curve reached 1.05 and 1.34, respectively, with the latter being significantly more efficient. This study achieved significant results in audio style conversion through the application of AutoML and big data analysis, successfully improving conversion accuracy. This progress has practical application value in multiple fields, including music production and sound effect design.

*Keywords*—*AutoML; audio style conversion; machine learning; big data analysis; adain module*

## I. INTRODUCTION

Audio style conversion, as an important branch in the field of audio processing, has always been the focus of researchers. The transformation of audio style aims to make subtle adjustments to the characteristics of audio without loss, such as time domain, frequency domain, timbre, pitch, etc., while retaining the essential information of audio [1-2]. The implementation of this transformation has a profound impact on many fields such as music production, speech synthesis, and oral teaching [3-4]. However, traditional audio style conversion methods often require a large amount of manual feature extraction and complex algorithm design. This limits the research process of audio style conversion. At present, to solve the above problems, most scholars advocate the introduction of data analysis to achieve various audio processing. But in this way, automatic search and optimization of audio conversion models and parameters can be achieved [5-6]. Meanwhile, this study also extracts valuable style information from massive audio data through big data analysis, further improving the accuracy and naturalness of audio style conversion. Integrating AutoML into audio style conversion research directly addresses the inefficiencies of current methodologies. This provides a systematic approach to model selection and parameter tuning, which is critical for enhancing the practicality and accessibility of audio style transformations. The innovation of research is mainly manifested in two aspects. AutoML is introduced into the research of audio style conversion to achieve automated search and optimization of audio conversion models, with the aim of improving the efficiency of audio style conversion. The second is to use big

data analysis methods to extract style information from massive audio data, making audio style conversion more accurate and natural. The research contributions consist of developing an AutoML-based framework for optimizing audio style conversion models efficiently, introducing big data analytics for extracting precise style features, establishing a benchmark dataset for comparative analysis that demonstrates enhanced conversion accuracy and naturalness, and validating real-world applications across music production, speech synthesis, and language learning. This study also provides new research methods and ideas for other related fields, with broad application prospects and important academic value. The research will be conducted in four parts. The first part is an overview of audio style conversion on the grounds of AutoML and big data analysis. The second part is the research on audio style conversion on the grounds of AutoML and big data analysis. The third is the experimental verification of the second. The fourth is a summary of the research content and points out the demerits.

## II. RELATED WORKS

Audio style conversion has always been an essential research topic in the audio processing, with the goal of achieving lossless conversion of audio styles while maintaining the original audio information. Li J et al. presented a novel ALRW method. The research results indicate that this method could markedly decrease compensation information. And it exhibits strong robustness to common operations. In the absence of an attack, it is possible to recover the covered audio signal without loss [7]. Lin F et al. presented a new text audio sentiment analysis framework called StyleBERT, which enhances unimodal sentiment information representation by learning different modal styles and reduces dependence on fusion. The research results indicate that StyleBERT performs excellently on multiple benchmark datasets, markedly superior to state-of-the-art multimodal baselines, and is an effective multimodal sentiment analysis framework [8]. Chen B and other scholars proposed a non-parallel data to speech conversion technology on the grounds of data augmentation - ParaGen. The experiment showcases that ParaGen can effectively convert the speaker identity of the source speech to the target speaker while preserving the local speaking style. And the converted speech possesses more excellent speech naturalness and speaker similarity than the statistical parameter speech synthesis system [9]. Xu D et al. proposed a bipolar phase shift modulation single-stage inverter for efficient and low distortion audio amplification. The research results were validated through a prototype with an output power of 200kHz and 250W. It demonstrates the effectiveness of the proposed

BPSM • FBAC-SSI method in improving the efficiency of audio amplifiers and reducing distortion [10]. Chandrakar R et al. proposed an enhanced system for automatic motion object detection and tracking using RBF-FDLNN and CFR algorithms. It can effectively handle the problem of motion target detection and tracking in traffic monitoring. The research results indicate that the proposed RBF-FDLNN classifier performs better than other existing methods in video frame object detection, proving the effectiveness of this method [11].

However, traditional audio style conversion methods rely on complex algorithm design and a large amount of manual feature extraction. This to some extent limits the development of audio style conversion technology. Zhang J presented a music feature extraction and classification model on the grounds of convolutional neural networks. The research results indicate that this method outperforms traditional manual models and machine learning based methods in music feature extraction and classification. This effectively addresses the shortcomings of traditional methods in feature selection and multi classification [12]. Singh P K et al. proposed new feature descriptor-binary image symbolization, for recognizing handwritten digits of different texts. The research results indicate that the symbolic feature descriptors of binary images have high script invariance, and can maintain high recognition rates even in mixed use of text [13]. Jiang ZG et al. proposed a segmentation and keyframe extraction method for video behavior recognition, and further proposed an improved vehicle detection algorithm on the grounds of fast R-CNN. The research results indicate that the application of keyframe extraction technology and optimized fast R-CNN model significantly improves the accuracy of vehicle detection, reduces missed detections, and demonstrates satisfactory detection rates [14]. Jia Z et al. proposed domain invariant feature extraction and fusion. The research results indicate that domain invariant feature extraction and fusion methods have achieved significant performance improvements on multiple datasets, effectively addressing the challenge of cross domain character re recognition [15]. Grzegorowski et al. proposed a supply management solution that considers individual delivery plans for each location. The research outcomes demonstrate that the method could markedly handle high uncertainty in data and effectively solve the cold start problem of vending machine networks [16].

Wu SL et al. utilized the advantages of Transformer and VAE to propose MuseMorphose for music generation, which is characterized by the user's ability to control style attributes. The results showed that MuseMorphose exceeded the RNN baseline in style transfer metrics [17]. Rashid A B et al. proposed an automatic detection model for student learning style in a learning management system based on online learning activities. The research shows that this model can assist educators in optimizing teaching content and recommending suitable learning materials based on student characteristics [18]. Chen et al. proposed reinforcement learning based audiovisual speech recognition framework MSRL, which focuses on stable supplementary information of visual modalities. The research results show that MSRL achieves the best performance on the LRS3 dataset, especially

demonstrating better universality in unknown noise testing [19].

In summary, existing research results indicate that AutoML can achieve automatic recognition and conversion of audio styles, which helps to solve the efficiency and accuracy problems of traditional methods in large-scale data processing. However, the complexity and diversity of audio data processing, such as feature extraction and model selection, remain challenges that limit the comprehensive application of AutoML. These technologies also need to be further optimized in practical scenarios such as music creation and speech synthesis. In view of this, this study aims to develop stronger audio feature extraction algorithms and establish effective model evaluation methods. And it is necessary to study how to better integrate these technologies into practical applications to maximize the potential of AutoML in audio processing. AutoML's advancement in audio style conversion heralds new creative horizons in music production, elevates speech synthesis realism, and promises tailored, immersive language learning experiences that are revolutionizing multiple industries.

## III. AUDIO STYLE CONVERSION METHODS ON THE GROUNDS OF AUTOML AND BIG DATA ANALYSIS

This study combines an improved VGG and EfficientNet feature extraction network to deeply extract audio data features. It utilizes Adain based normalization modules and feature decoding networks to achieve lossless audio style conversion. It combines AutoML and big data analysis to construct an automatically optimized audio style conversion model to improve conversion efficiency and accuracy. This study integrates the latest machine learning techniques into audio processing, providing a new research perspective for the development of audio style conversion technology.

### A. Based on Improved VGG and EfficientNet Feature Extraction Network

Audio style conversion relies on deep learning to automatically extract features, obtaining abstract and robust features. The VGG network has fewer parameters and requires less computing resources. The new EfficientNet breaks the convention of improving network performance in a single dimension by adjusting input resolution, depth, and width, achieving a balance between accuracy and efficiency [20-21]. The VGG-16 network uses small convolutional kernels instead of large ones to enhance model nonlinearity, reduce computational complexity, and remove fully connected layers (FCL). Then it changes the pooling layer to a convolutional layer with a stride of 2, and uses a swish activation function to improve model performance. EfficientNet extracts useful features in audio style conversion, compares feature representations, and predicts the effect of style conversion. The optimization process relies on the loss function, and the smaller the loss function, the higher the accuracy of the model. Therefore, EfficientNet and VGG networks have important application value in the study of audio style conversion [22-23]. The mean square error is shown in Eq. (1).

$$MSE = \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - f\left(x_i\right) \right)^2 \qquad (1)$$

In Eq. (1), $y_i$ serves as the actual value, $f(x_i)$ serves as the predicted value, and $f(x_i)$ serves as the number of samples. The backpropagation of gradient information is crucial for neural network algorithms to self-learn and update. The optimized EfficientNet algorithm is showcased in Eq. (2).

$$\theta_k = \theta_{k-1} - \alpha \cdot g \tag{2}$$

In Eq. (2), $\theta_k$ is the parameter value at the current time, $\alpha$ is the learning rate, and $g$ is the gradient. It increases the number of audio processing channels and adds feature layers to obtain more audio features. It increases network depth, utilizes deep neural networks to improve performance, and enhances audio feature extraction. It improves the input audio sampling rate, enhances network accuracy, enriches audio features, and reduces information loss. From this, it can be concluded that the tensor of the network output is shown in Eq. (3).

$$Y_i = F_i(X_i) \tag{3}$$

In Eq. (3), $X_i$ is the tensor of a specific convolutional layer. A deep neural network composed of $k$ convolutional layers is shown in (4).

$$N = F_k \odot \ldots \odot F_2 \odot F_1(X_1) = \underset{j=1\ldots k}{\odot} F_j(X_1) \tag{4}$$

In Eq. (4), $\odot$ is the multiplication operation, $i$ is the stage number, and $i$ is a single operation. Scaling the model could enhance the accuracy of the network within the limits of memory and computational complexity, as shown in Eq. (5).

$$\begin{cases} \underset{d,w,r}{\max} \ Accuracy\big(N(d,w,r)\big) \\ s,t, N(d,w,r) = \underset{i=1\ldots S}{\odot} F_i^{d,\hat{L}}\Big(X_{r \cdot \hat{H}_i^r \hat{W}_i^R \hat{C}_i}\Big) \end{cases} \tag{5}$$

In Eq. (5), $d$, $w$, and $r$ represent the scaled depth, width, and resolution, respectively. In EfficientNet, achieving composite optimization by simultaneously scaling three dimensions at appropriate proportions could enhance the performance and classification accuracy of the network. This could also decrease the computational complexity of the model and enhance the performance of the network. The MBConv module used internally is the core structure, which is a unique feature extraction structure of EfficientNet. The two-dimensional view is efficiently extracted during the continuous stacking process in the Block layer. The MBConv module is shown in Fig. 1.

In Fig. 1, EfficientNet first performs pointwise convolution on the input feature map and adjusts the expansion ratio by changing the output channel dimension. Then it performs deep convolution, reducing the dimension to the original number of channels, and then performs point by point convolution again. This network module integrates compression and arousal of network attention to focus on channel features. The feature map is processed by stacking 32 MBConvs, and then sequentially passes through one-dimensional convolutional layers, global average pooling 2D, and FCL to generate feature vectors with a dimension of 2640. EfficientNet reduces the computational complexity of the network through deep convolution and point by point convolution, compared to conventional convolution operations. The schematic diagram of EfficientNet network structure is showcased in Fig. 2.



Fig. 1. MBConv module.

Fig. 2.    EfficientNet network structure diagram.

In Fig. 2, EfficientNet utilizes MBConv as the backbone network, which originates from MobileNet V2. MBConv includes a regular convolution, a deep convolution (including BN and Swish), an SE module, another regular convolution (for dimensionality reduction, including BN), and a Droupout layer. The SE module contains a global average pooling and two FCL. The quantity of nodes in the first FCL is equal to the quantity of channels in the feature matrix of the input MBConv, and the activation function is Swish. The quantity of nodes in the second FCL is equal to the number of channels in the output feature matrix of the deep convolutional layer, and the activation function is sigmoid.

### B. Audio Style Conversion Normalization Module and Feature Decoding Network on the Grounds of Adain

After completing the feature extraction for audio style conversion, the next step is to use the Adam based normalization module and feature decoding network for audio style conversion. In this process, Adain technology is used to convert audio features into styles, and then a feature decoding network is used to convert the converted features into perceptible audio signals. This can achieve style conversion of audio.

The normalization process can make the data distribution have a mean of 0 and a variance of 1, which helps to avoid gradient vanishing and exploding, thus accelerating the training process. When processing large amounts of data, BatchNorm needs to use mini batch data to estimate mean and variance, but training may become unstable when computing power is limited and the input audio data volume is too large. The Adain method confirms that the Instance Normalization layer can reduce style loss faster than the BN layer, thereby accelerating training. The core of Adain is to fuse the features obtained from content audio and style audio through an Encoder network, and then decode them to obtain style audio. The decoding of style audio is shown in Eq. (6).

$$AdaIN(x, y) = \sigma(y)\left(\frac{x - \mu(x)}{\sigma(x)}\right) + \mu(y) \tag{6}$$

In Eq. (6), the content image is $x$ and the style image is $y$. The current mainstream normalization methods mainly include Batch Normalization, Layer Normalization, Instance Normalization, Group Normalization, and Switchable Normalization. These methods are all on the grounds of normalization processing of different dimensions of input audio. Specifically, given the dimensions of the input audio as (N, C, H, W), different normalization methods choose different normalization strategies on these four dimensions. An example of centralized normalization is shown in Fig. 3.



Fig. 3.    An example of centralized normalization diagram.

Batch Normalization is the normalization of NHW on each batch. Due to its calculation of mean and variance on each batch, if the batch size is too small, the calculated mean and variance may not represent the distribution of the entire data, which may lead to unstable training and poor performance. Layer Normalization is the normalization of CHW for each channel direction, mainly applied in RNN networks. Compared to Batch Normalization, Layer Normalization solves the problem of deep non fixed networks by normalizing all neurons in each layer of the deep network, as shown in Eq. (7).

$$\begin{cases} \mu^l = \dfrac{1}{H} \sum_{i=1}^{H} a_1^l \\ \sigma^l = \sqrt{\dfrac{1}{H} \sum_{i=1}^{H} \left( a_i^l - \mu^l \right)^2} \end{cases} \quad (7)$$

In Eq. (7), $\mu$ is the mean and $\sigma^l$ is the variance. Instance Normalization is mainly applied in the field of audio style conversion, which normalizes audio signals at the pixel level. Due to the fact that the generated results mainly rely on specific audio samples, it is very suitable for audio stylization, which can accelerate model convergence and maintain independence between samples. Group normalization is achieved by grouping channels and normalizing them within the group, and its calculated mean is independent of batch size. Therefore, it can solve the impact of batch normalization on training results when the batch is small. In the feature map of each sample, channels are divided into G groups, each containing C/G channels, and the mean and standard deviation of these channels are calculated. Switchable Normalization combines BN, LN, and IN normalization methods, assigning them different weights to enable the network to learn which normalization method to use on its own, thus achieving adaptive selection of normalization methods. The style transition effect is showcased in Fig. 4.

The decoder and encoder have a symmetrical structure. During the audio style conversion, the encoder is responsible for feature extraction of the original audio and the target style audio. The decoder combines the original audio features and style features generated by Adain to generate stylized audio. In audio style conversion systems, only the decoder is usually trained, while the parameters of the feature extraction and loss calculation networks remain unchanged. During the feature extraction, downsampling is usually performed through a series of convolutions. In the feature decoding stage, it is necessary to upsample the features to restore the size of the original audio. Common upsampling methods include linear interpolation, deconvolution, and deconvolution. Deconvolution is a special type of convolution that fills feature audio with zeros and then convolves it by rotating the convolution kernel. In decoding networks for audio style conversion, interpolation algorithms are commonly used for upsampling operations. Deconvolution restores feature audio, which operates in the opposite direction of convolution and is essentially transposed convolution. The relevant schematic diagram is showcased in Fig. 5.

### C. Audio Style Conversion Model on the Grounds of AutoML and Big Data Analysis

Considering the characteristics of audio data and the complexity of processing audio data, this study selected an audio style conversion model on the grounds of AutoML and big data analysis for research. The model utilizes big data analysis technology to process massive audio data. And it automatically finds the optimal model parameters through AutoML to achieve more efficient and accurate audio style conversion. The framework structure of the model is set as showcased in Fig. 6.



Fig. 4. Style transition effect diagram.



Fig. 5. The relevant schematic diagram.



Fig. 6. Functional module execution process.

In Fig. 6, considering the characteristics of audio style conversion, the audio style conversion process on the grounds of AutoML and big data analysis can be mainly divided into four stages: preprocessing stage, feature extraction stage, model training stage, and style application stage. The preprocessing stage is mainly used for cleaning and formatting audio data for subsequent processing. The feature extraction stage converts audio data into feature vectors that can be processed by machine learning models. During the model training phase, AutoML is used to automatically search for the optimal model parameters, to achieve more efficient and accurate audio style conversion. The final style application stage applies the trained model to new audio data to complete style conversion. The evaluation indicators for audio style conversion are shown in Eq. (8).

$$D_i = \frac{m_i}{N_i} \tag{8}$$

In Eq. (8), $m_i$ represents the degree of style distortion of the audio sample, and $N_i$ represents the total number of style transitions performed. It enhances the accuracy of the model by connecting all the encoded features obtained, and the vector construction is shown in Eq. (9).

$$\begin{cases} G(q_t, a_t) = q + (\max(q) + 1) \cdot a_t \\ v_t = Q(C(q_t, a_t)) \cap Q(C(f_t, a_t)) \cap Q(dc_t) \end{cases} \tag{9}$$

In Eq. (9), $q_t$ represents the feature number, and $dc_t$ represents the difficulty coefficient of coherent instructions. $G(\ )$ represents instruction execution combination, $Q(\ )$ represents encoding format, and $\cap$ represents connection. It uses ReLU activation function to train the autoencoder and removes the output layer after completion. Then it uses the output of the hidden layer as input in the AutoML model, as shown in Eq. (10).

$$v_t^{'} = FAKT(W \cdot v_t + b) \tag{10}$$

In Eq. (10), $W, b$ represent the weight matrix and bias vector for audio execution, respectively. The current state is obtained by combining the information processed by the forget gate with the input information obtained by the input gate. The process is shown in Eq. (11).

$$\begin{cases} f_t = \sigma \left( W_f \left[ v_t^{'}, h_{t-1} \right] + b_f \right) \\ c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \end{cases} \tag{11}$$

In Eq. (11), the output gate $o_t$ determines what information to extract from $c_t$ on the grounds of $h_{t-1}$, $v_t^{'}$, and the ReLU function, forming a hidden state $h_t$. Then, combined with historical encoding, the attention score is calculated using weight factors, as shown in Eq. (12).

$$s\left(v_i^{'}, v_t^{'}\right) = v^T \tanh\left(W_1 v_i^{'} + W_2 v_t^{'}\right) \tag{12}$$

In Eq. (12), $v_i^{'}$ represents the encoding of historical sequence information on the grounds of data compression, and $v_t^{'}$ represents the dimensionality reduction encoding of input information at time $t$. $W, v$ are both network parameters for the Department of Science. After weighting, cross entropy can be used in machine learning for measuring the difference between actual labels and predicted results. Therefore, the study uses it as the loss function of the AutoML model, as shown in Eq. (13).

$$L = -\sum_t \left( a_{t+1} \log y_t^T \delta(q_{t+1}) + (1 - a_{t+1}) \log\left(1 - y_t^T \delta(q_{t+1})\right) \right) \tag{13}$$

In Eq. (13), $y_t^T \delta(q_{t+1})$ and $a_{t+1}$ represent the predicted and true probability distributions, respectively. It assumes that $X$ is one knowledge unit, including $n$ execution nodes. And if the probability of using the execution node $x_i$ is $P(x_i), i = 1, 2, \ldots, n$, then the audio modeling's mastery of the execution node is shown in Eq. (14).

$$Z(X) = -\sum_{i=1}^n P(x_i) \log_2 \frac{A(x_i)}{2} \tag{14}$$

In Eq. (14), $A(x_i)$ represents the execution efficiency of the audio at the execution node. $P(x_i)$ represents the probability of the execution node appearing.

## IV. Modeling and Analysis of Audio Style Conversion on the Grounds of AutoML and Big Data Analysis

It conducts research on audio style transformation modeling on the grounds of AutoML and big data analysis, and extracts deep features from audio data. Then it utilizes EfficientNet and VGG networks to construct an audio style classification and transformation model. By comparing different audio features, it predicts the performance of audio after different style transitions. The experimental environment and dataset parameters are showcased in Table I.

The accuracy changes of the audio style conversion model on the conversion of training set audio and validation set audio under different iterations are shown in Fig. 7.

TABLE I.        EXPERIMENTAL ENVIRONMENT AND DATASET PARAMETERS

| Parameter | Description | Parameter | Description |
|---|---|---|---|
| Operating System | Ubuntu 20.04 | CUDA Version | 11.2 |
| CuDNN Version | 8.1.0 | Programming Language | Python 3.8 |
| Parameter | Description | Parameter | Description |
| Operating System | Ubuntu 20.04 | CUDA Version | 11.2 |
| CuDNN Version | 8.1.0 | Programming Language | Python 3.8 |
| Data Augmentation | Add noise, Time stretch | Audio Sample Rate | 44100 Hz |
| Audio Resolution | 16 bits | Style Audios Styles | Classical, Rock, Jazz, Pop |
| Total Samples in Training Data | 1.2 million samples | Training Data Size after Processing | 1024*1024 |
| Processed Sample Count | 16 samples per audio | Scales during Training | 400400, 300300, 256*256 |
| Style Audio Size | 512*512 | Big Data Analysis Tool | Apache Hadoop, Apache Spark |



Fig. 7.    Accuracy of training and validation sets in audio style transformation.

In Fig. 7, with the quantity of iterations grows, the accuracy of training audio style conversion gradually increases. This indicates that the audio style conversion model has a faster convergence ability for training audio and verifying audio. However, the accuracy of verifying audio style conversion fluctuates greatly. Sometimes the conversion accuracy of the verified audio is higher than that of the training audio, but lower than that of the training audio at other iterations. The fluctuation amplitude gradually decreases with the increase of iterations. After 50 iterations of training, the style conversion accuracy of both the training audio and validation audio exceeded 90%, and the effect was significant. The training curve tends to stabilize after 100 iterations, while the validation curve reaches stability after 175 iterations. These two curves indicate that the audio style conversion model achieved good training performance after 175 iterations, and reached its optimal performance after 200 iterations. For the audio samples of bird singing, vehicle horn sound, wave crashing sound, and piano performance, the style conversion ratio is compared with the original style conversion framework, as shown in Fig. 8.

In Fig. 8, the curves of style conversion ratios are higher than those of the original framework, demonstrating that the style conversion model could improve the conversion quality of audio. This is because the model fully takes into account the characteristics of audio style transition, that is, in the audio, certain parts of the style transition are more pronounced than others. It calculates the weighted value of each audio segment on a unit basis. Then it adaptively compensates for the weighted values of each segment, making the style transformation of each segment's weighted values more detailed. This is to achieve the goal of improving the quality of audio style conversion and enhancing audio performance. It compares the improvement in style conversion efficiency for three audio sample sets: FreeSound, Looperman, and SampleSwap, as shown in Fig. 9.



Fig. 8.    Comparison of style transformation ratio with original framework.

In Fig. 9, the yellow curve represents the efficiency of SampleSwap style conversion. The blue dashed line represents the efficiency of Loopperman style conversion. The green curve represents the efficiency of FreeSound style conversion. In the efficiency analysis curve of Fig. 11 (a), the green curve has the best effect and is also the most stable, with time ranging from 0 to 300. The green curve is the fastest to reach 4.48 and has been operating at this efficiency. Compared to the yellow lines, the efficiency of the blue dashed lines is much lower. In the efficiency analysis curve of Fig. 11 (b), the yellow curve has the worst effect, only less than 60, while the green curve and blue dashed line are 140 and 116, respectively. In Fig. 11 (c), the efficiency of the yellow curve is around 1.05, while the efficiency of the green curve is still as high as 1.34. The spectrograms of the flute version of "Butterfly Lovers", the violin version of "My Heart Forever", and the flute version of "Titanic" output by the STFT model are shown in Fig. 10.

(a) Efficiency Analysis in Audio
sample bird song



(b) Efficiency Analysis in
Vehicle horn sound



(c) Effectiveness Analysis in the
sound of waves crashing on the shore

Fig. 9.  Efficiency analysis of audio style transformation improvement in three datasets.



(a) The flute version of 'Liang Zhu' output from
the STFT model



(b) The violin version of 'My Heart Forever' output from the STFT
model



(c) The violin version of Titanic output from the
STFT model

Fig. 10.  The spectrograms of the flute version of "Butterfly Lovers", the violin version of "My Heart Forever", and the flute version of "Titanic" output from the
STFT model.

In Fig. 10, whether it is the flute to violin or the violin to flute, the audio quality of the spectrograms output by the two models is not very good, resulting in poor sound quality. The audio obtained by the STFT model hardly shows any changes in timbre, and the sound is relatively noisy. The audio obtained by the CQT model can vaguely distinguish the timbre of the instrument. The time domain diagram is drawn on the grounds of the first eight seconds of Beethoven's First String Trio in e-flat major (hereinafter referred to as Beethoven. wav) and the first nine seconds of Telemann's Flute Fantasy. The time-domain diagram drawn by Beethoven. wav and Telemann in the first nine seconds are shown in Fig. 11.

(a) Beethoven. wav Time Domain DiagramStress response curve of nodes at the point of maximum stress

(b) Telemann time-domain diagram

Fig. 11. Time-domain plots plotted from the first 9s fragments of beethoven.wav and Telemann.

In Fig. 11, the horizontal axis serves as time and the vertical axis serves as amplitude, reflecting the temporal variation of the audio signal. In this figure, the amplitude variation of Beethoven. wav is relatively stable, indicating the smooth and harmonious nature of Beethoven's trio. Telemann's amplitude changes significantly, showcasing the dynamic changes and rhythmic sense of flute fantasies.

## V. RESULTS AND DISCUSSION

The results of applying AutoML and big data analysis to audio style conversion are presented. The results show that AutoML can recognize and transform audio styles more efficiently than traditional methods. By automating the process of model selection and feature extraction, the method has greatly improved the efficiency and accuracy of processing large-scale data sets. The analysis of a significant amount of audio data shows that the feature extraction algorithm developed in this study is robust and capable of capturing the fine features required for high-fidelity audio style conversion. In addition, the established evaluation method has proven effective in determining the optimal model across different data sets and transformation tasks. Challenges remain, especially with regard to the complexity and diversity of processing audio data. Despite the progress made, the deployment of these technologies in real-world scenarios such as music creation and speech synthesis needs to be further optimized. It is evident that while AutoML simplifies workflows, the complex nature of audio data requires a sophisticated understanding of domain-specific functionality, which is not fully implemented by the current AutoML framework. The discussion highlighted the importance of improving these technologies for practical applications. AutoML's potential for audio processing is enormous, but its full application is limited by current technology. Future research should therefore focus on improving AutoML's adaptability to the specific needs of audio data, ensuring that the benefits of automation can be fully utilized in both practical and creative contexts. The integration of these advanced technologies will have a major impact on areas such as music production, speech synthesis, and more, as long as subtle adaptations are taken into account.

## VI. CONCLUSION AND FUTURE WORK

Audio style conversion is an important research field in digital audio processing, with the goal of changing the style characteristics of audio content without altering it. Audio style conversion on the grounds of AutoML and big data analysis can automatically learn and convert audio styles, thereby improving the efficiency and quality of audio processing. The research results show that using iterative learning for audio style conversion training, the training curve tends to stabilize after 100 iterations, while the validation curve reaches stability after 175 iterations. In efficiency analysis, the efficiency of the yellow curve and the green curve reached 1.05 and 1.34, respectively, with the latter having significantly higher efficiency. In the audio analysis section, some parts had more obvious style transitions than others, such as Telemann's significant amplitude changes, showcasing the dynamic changes and rhythm of flute fantasies. The main contribution of the research lies in utilizing AutoML and big data analysis methods for enhancing the accuracy and efficiency of audio style conversion, offering new tools and methods for music production and sound effect design. However, this study also has some shortcomings, such as poor style switching effects in certain parts of the audio, and the need to improve sound quality. There is still a lot of room for advancement in the study of audio style conversion in future research. It is necessary to further optimize and improve the methods of AutoML and big data analysis to enhance the accuracy and efficiency of audio style conversion. It also brings new possibilities for exploring the application of audio style conversion in more fields, such as speech synthesis, music generation, entertainment industry, etc.

## REFERENCES

[1] C. A. Hallin, M. Koren, A. A. Issa, O. Koren, "AutoML classifier clustering procedure," International Journal of Intelligent Systems, vol. 37, pp. 4214-4232, 2022.

[2] H. Cai, J. Lin, Y. Lin, Z. Liu, H. Tang, H. Wang, L. Zhu, S. Han, "Enable Deep Learning on Mobile Devices: Methods, Systems, and Applications," ACM Transactions on Design Automation of Electronic Systems, vol. 27, pp. 20-50, 2022.

[3] H. S. Yang, K. R. Kim, S. Kim, J. Y. Park, "Deep Learning Application in Spinal Implant Identification," Spine, vol. 46, pp. 318-324, 2021.

[4] O. Owoyele, P. Pal, A. V. Torreira, D. Probst, M. Shaxted, M. Wilde, P. K, "Senecal. Application of an automated machine learning-genetic algorithm (AutoML-GA) coupled with computational fluid dynamics simulations for rapid engine design optimization," International Journal of Engine Research, vol. 23, pp. 1586-1601, 2021.

[5] M. K. Shende, A. E. Feijoo-Lorenzo, N. D. Bokde, "cleanTS: Automated (AutoML) Tool to Clean Univariate Time Series at Microscales," Neurocomputing, vol. 500, pp. 155-176, 2021.

[6] X. He, K. Zhao, X. Chu, "AutoML: A survey of the state-of-the-art," Knowledge-Based Systems, vol. 212, pp. 1-27, 2021.

[7] J. Li, S. Xiang, "Audio-lossless robust watermarking against desynchronization attacks," Signal Processing, vol. 198, pp. 108561-108573, 2022.

[8] F. Lin, S. Liu, C. Zhang, J. Fan, Z. Wu, "StyleBERT: Text-audio sentiment analysis with Bi-directional Style Enhancement," Information systems, vol. 114, pp. 1-11, 2023.

[9] B. Chen, Z. Xu, K. Yu, "Data augmentation based non-parallel voice conversion with frame-level speaker disentangler," Speech Communication, vol. 136, pp. 14-22, 2022.

[10] D. Xu, S. Zhong, J. Xu, "Bipolar Phase Shift Modulation Single-Stage Audio Amplifier Employing a Full Bridge Active Clamp for High Efficiency Low Distortion," IEEE Transactions on Industrial Electronics, vol. 68, pp. 1118-1129, 2021.

[11] R. Chandrakar, R. Raja, R. Miri, U. Sinha, A. K. S. Kushwaha, H. Raja, "Enhanced the moving object detection and object tracking for traffic surveillance using RBF-FDLNN and CBF algorithm," Expert Systems with Applications, vol. 191, pp. 1-15, 2022.

[12] J. Zhang, "Music Feature Extraction and Classification Algorithm Based on Deep Learning," Scientific Programming, ,vol. 2, pp. 1-9, 2021.

[13] P. K. Singh, I. Chatterjee, R. Sarkar, E. B. Smith, M. Nasipuri, "A new feature extraction approach for script invariant handwritten numeral recognition," Expert Systems, vol. 38, pp. 1-22, 2021.

[14] Z. G. Jiang, X. T. Shi, "Application Research of Key Frames Extraction Technology Combined with Optimized Faster R-CNN Algorithm in Traffic Video Analysis," Complexity, vol. 4, pp. 1-11, 2021.

[15] Z. Jia, Y. Li, Z. Tan, W. Wang, Z. Wang, G. Yin, "Domain-invariant feature extraction and fusion for cross-domain person re-identification," The visual computer, vol. 39, pp. 1205-1216, 2023.

[16] M. Grzegorowski, J. Litwin, M. Wnuk, M. Pabi, U. Marcinowski, "Survival-Based Feature Extraction—Application in Supply Management for Dispersed Vending Machines," IEEE transactions on industrial informatics, vol. 19, pp. 3331-3340, 2023.

[17] S. L. Wu, Y. H. Yang, "MuseMorphose: Full-song and fine-grained piano music style transfer with one transformer VAE," IEEE/ACM Trans. on Audio, Speech, and Language Processing, vol. 31, no. 5, pp. 1953-1967, 2023.

[18] A. B. Rashid, R. R. R. Ikram, Y. Thamilarasan, L. Salahuddin, N. F. Abd Yusof, Z. B. Rashid, "A Student Learning Style Auto-Detection Model in a Learning Management System," Eng. Tech. & Appl. Sci. Res., vol. 13, no. 3, pp. 11000-11005, 2023.

[19] C. Chen, Y. Hu, Q. Zhang, H. Zou, B. Zhu, E. S. Chng, "Leveraging modality-specific representations for audio-visual speech recognition via reinforcement learning," Proc. of the AAAI Conf. on Artificial Intelligence, vol. 37, no. 11, pp. 12607-12615, 2023.

[20] A. K. Nair, J. Sahoo, E. D. Raj, "Privacy preserving Federated Learning framework for IoMT based big data analysis using edge computing," Computer Standards and Interfaces, vol. 86, pp. 1-20, 2023.

[21] B. Wang, J. Wan, Y. Zhu, Y. Chen, "Institutional capability, cooperation level and irrigation water order: Empirical analysis based on survey data from the Yellow River area," Irrigation and Drainage, vol. 72, pp. 716-728, 2023.

[22] Y. Fang, B. Luo, T. Zhao, D. He, B. Jiang, Q. Liu, "ST-SIGMA: Spatio-temporal semantics and interaction graph aggregation for multi-agent perception and trajectory forecasting," CAAI Transactions on Intelligence Technology, vol. 7, pp. 744-757, 2022.

[23] J. Purohit, R. Dave, "Leveraging Deep Learning Techniques to Obtain Efficacious Segmentation Results," Archives of Advanced Engineering Science, vol. 1, pp. 1-16, 2023.

# Attraction Recommendation and Itinerary Planning for Smart Rural Tourism Based on Regional Segmentation

Ruiping Chen[1*], Yanli Zhou[2], Dejun Zhang[3]

Compliance Office, Foshan Polytechnic, Foshan, 528137, China[1]

School of Culture Tourism and Creativity, Foshan Polytechnic, Foshan, 528137, China[2, 3]

*Abstract*—As the rural tourism industry develops, effective attraction recommendations and planning are crucial for the tourist experience. Then, a rural scenic spot tourism recommendation and planning technology based on regional segmentation was proposed. The scenic area was divided into multiple grids based on tourist check-in behaviour, and the interest and influence of the scenic area were associated with the grid check-in behaviour. Content recommendation was achieved through two factors: popularity and regional location. And considering the sparsity of data in the recommendation, clustering algorithms were introduced to model tourist check-in behaviour based on factors such as time and regional location, and content recommendation was achieved through tourist preferences. In the performance analysis of recommendation models, the proposed model has an accuracy of 0.965 and 0.956 on the Gowalla and Yelp datasets, respectively, which is superior to other models. Comparing the recommendation loss performance of different models, the proposed model has an RMSE loss of 0.120 on the Gowalla dataset, which is superior to other models. In practical application analysis, when the recommended number is 5, the accuracy and recall of the proposed model are 0.138 and 0.069, respectively, which are superior to other models. In tourism itinerary planning, the overall planning time of the model is the shortest. Therefore, the proposed model has excellent application effects, and the research content provides important technical references for tourist travel and rural tourism destination planning.

*Keywords—Regional division; trip planning; recommended tourist attractions; clustering algorithm; time factor*

## I. INTRODUCTION

According to data released by the National Tourism Administration, the average annual growth of rural tourism tourists has exceeded 10%, becoming an important industry supporting local economic development. However, traditional recommendation techniques mainly rely on user historical preferences and ratings to make recommendations, ignoring the changing interests of tourists at different times and in different regions [1]. In addition, due to the relatively small amount of data on rural tourist attractions and the sparsity of data, traditional recommendation systems couldn't meet the tourism needs of tourists [2]. The existing tourism recommendation methods may ignore user preferences, and their recommendation content may be inaccurate. To address this issue, Tourist Sign-in Area Segmentation (TSAS) has been proposed. This technology uses the check-in data of

tourists to divide the scenic area into multiple grids and associates the interest and influence of the scenic area by analyzing the check-in behaviour of tourists [3]. In addition, to address data sparsity, this technology introduces clustering algorithms to model the check-in behaviour of tourists and recommends content based on their preferences, improving the recommendation accuracy. There are several innovations in the recommendation technology studied. Firstly, it combines geographical location and time factors to achieve more accurate recommendations and planning of rural tourism attractions. Through in-depth analysis of tourist check-in behaviour, this technology can accurately capture the interests and preferences of tourists and provide personalized recommendations and planning solutions based on factors such as time and geographical location. The research content will provide reference for recommending rural tourism attractions and tourist itinerary planning and accelerate the development of the rural tourism industry.

The research content includes four sections. Introduction is given in Section I, Related works is given in Section II. The construction of rural tourism recommendation and itinerary planning model based on regional segmentation is given in Section III. Section IV apply the mentioned technology to specific scenarios and verify the effectiveness of the proposed recommendation model in practical scenarios. Finally, Section V gives the summary and analyses of the entire article are conducted, and the direction of technological improvement in the future is elaborated.

## II. RELATED WORKS

Recommendation technology is mainly used to solve information data overload, which can help users find suitable information content faster and more accurately. At present, recommendation technology has been widely applied in various fields, and relevant scholars have conducted extensive research on it. Cui Z et al. found that traditional recommendation systems may overlook the inherent relationship between user preferences and time. To address this problem, a new fusion recommendation model based on time correlation coefficients was proposed. This model further improved the accuracy and efficiency of recommendations by clustering similar users together. In addition, the study also proposed a personalized recommendation model based on preference patterns, mainly analyzing user behaviour to optimize content recommendation. The effectiveness of the

proposed model was validated using two datasets, MovieLens and Douban. Compared to other models, the overall recommendation performance of the proposed model was better [4]. Zhou X et al. focused on modeling and analyzing patient doctor generated data using an ensemble-based deep learning framework. So a fusion extraction model was proposed in the study, which could extract and highlight semantic information in patient inquiries. Then, the study proposed an intelligent recommendation method that refined the learning process through clustering mechanisms, providing patients with automatic clinical guidance and diagnostic recommendations. The accuracy of online patient queries could be effectively improved by applying the proposed technology to specific scenarios [5]. Cho J et al. focused on the impact of recommendation algorithms on user opinions on the video sharing platform YouTube. Therefore, traditional recommendation models were improved by processing information based on experimental results and filtering un-healthy content information. Through testing, the researched technology had better recommendation performance in practical scenarios and could filter out the impact of harmful information on users [6].

Interest-based recommendation technology has a wide range of applications in tourism services and other fields. Interest-based recommendation technology focuses more on factors such as user preferences and behavioural habits, which is closer to the actual needs of users. Nitu P et al. conducted research on tourism recommendation technology based on social media activities. To improve the recommendation effect, personalized recommendations of the model were achieved by analyzing user Twitter data as well as research friend and follower data to identify travel-related tweets. Time-sensitive nearest degree weights were introduced to improve recommendation accuracy. The proposed technology applied to practical tourism recommendation scenarios had excellent recommendation performance, which was superior to other recommendation techniques [7]. Giglio S et al. conducted research on urban tourism recommendation technology. Clustering analysis was used to collect and analyze image data from multiple cities in Italy to improve the recommendation accuracy of the model, and Wolfram Mathematica was used to automatically identify clusters around points of interest. New tourism scenarios and more information for the interest point recommendation process could be provided by applying this technology to tourism recommendation scenarios, which was superior to other recommendation models [8]. Huang F et al. found that existing tourism route planning methods were mainly targeted at specific tasks and couldn't be applied to other tasks. To address this issue, a multi-task deep travel route planning framework was proposed, which integrated rich auxiliary information to construct a flexible planning model. These results confirmed that this method exhibited flexibility and superiority in travel route planning, outperforming relevant recommendation models [9]. Wang et al. focused on the importance of location in recommendation systems. The study proposed a multi-objective recommendation framework based on location and preference perception, modeling location-based recommendations as a multi-objective optimization problem. The study considered the performance of recommendation algorithms in recommending similar and different items, and a new multi-objective evolutionary algorithm was proposed. These results confirmed that this model could generate better recommendation solutions and overcome data sparsity and cold start issues compared to other recommendation models, resulting in better overall recommendation performance [10].

In summary, recommendation technologies mainly analyze the feature associations between objects and targets to achieve effective content recommendation. The above studies have analyzed the application of recommendation technology in different fields and discussed its effectiveness based on interest points. However, existing research has problems such as neglecting changes in user dynamic interests, insufficient understanding of user behaviour at a deeper level, and insufficient explanation of recommendation systems. Therefore, a tourism recommendation technology based on regional segmentation is proposed to provide important technical support for the development of the tourism industry and the promotion of tourism destinations.

## III. CONSTRUCTION OF RURAL TOURISM RECOMMENDATION AND ITINERARY PLANNING MODEL BASED ON REGIONAL SEGMENTATION

This section proposes a recommendation technique based on regional segmentation to segment rural areas and establish a recommendation model based on attendance. Simultaneously considering factors such as data sparsity and tourist interest transfer, the recommendation technology is improved and modeled based on time characteristics.

### A. Construction of a Recommendation Model Based on Tourist Check-in Area Segmentation

In recent years, the rural tourism industry has flourished, with a large number of tourists flocking to rural tourist attractions, promoting the development of the rural economy. To meet the personalized tourism needs of tourists, accurate recommendation of tourist attractions is crucial for improving the quality of travel. A recommendation method based on the division of tourist check-in areas has been proposed [11]. This method mainly considers that tourists will sign in and clock in when visiting the scenic area, share on their respective social circles or social circles, and use Location-based Social Network (LBSN) data information to mine tourist behaviour data. Segmentation is carried out according to the check-in area, and multiple areas are delineated to achieve content recommendation based on the size of regional influence [12]. The proposed TSAS method can accurately capture the interests and preferences of tourists by conducting in-depth analysis of their check-in behaviour in different regions. Unlike traditional recommendation systems, TSAS comprehensively considers geographical location and time, enabling recommendation systems to provide more timely and regional recommendations based on the geographical location and different periods of tourists, enhancing the personalization of recommendations. The relatively small and sparse data for rural tourism attractions can be addressed by using clustering algorithms to model the check-in behaviour of tourists, and the efficiency of recommendation systems in utilizing limited data is improved. Fig. 1 shows the framework of the entire tourism recommendation system.

Fig. 1 shows the framework of the entire tourism recommendation system, which implements content recommendation by mining feature data of tourists and rural tourism attractions and ranking them based on personalized feature influence. The preferred method is to obtain check-in information for rural tourism areas from LBSN data and segment the area based on the dimensions of the check-in area to obtain multiple small grid areas [13]. Fig. 2 is a schematic diagram of segmenting rural scenic spots.

According to Fig. 2, each small grid area contains the check-in information of tourists, which gathers various check-in interest points, and the characteristics of interest points between different grids are not the same. The length and width of the entire rectangular area are defined as $a$ and $b$. Two independent matrices need to be constructed after dividing the rectangular area into multiple grids, namely the tourist activity area matrix $X$ and the interest point area influence matrix $Y$. The matrix for a tourist $u$ in the activity area is denoted as $x_u$. For some tourists who have checked in the grid, the probability of tourists appearing in the area will be greater than 0 [14]. The influence vector of a certain interest point $l$ in the region matrix is set to $y_l$. The regional influence of scenic spots is mainly influenced by two factors: distance from surrounding locations and points of interest. The influence of interest point $i$ on the network region $l$ is represented by Eq. (1).

$$w_{ii} = p(i) \cdot \frac{1}{\sigma} K(\frac{d(i,l)}{\sigma}) \tag{1}$$



Fig. 1. Tourism recommendation system framework



Fig. 2. Schematic diagram of rural tourist attractions segmentation.

In Eq. (1), $d(i,l)$ is the distance from $i$ to the center of the grid. $p(i)$ represents the popularity of interest points. $\sigma$ means the standard deviation. $K(.)$ is a normal distribution. The number of tourist check-in is used as the popularity of the region, and the check-in data are normalized using Eq. (2).

$$p(i) = 1 + \lg(1 + r_{ui} \times 10) \tag{2}$$

In Eq. (2), $r_{ui}$ represents the cross-factor between geographic location and prevalence. Next, it is necessary to explore the relationship between the location and popularity of the region. The farther away the rural scenic spots are, the gradually decreasing influence can be considered. If they are too far away, the influence will be ignored. Taking the influence of interest points as a key consideration, the influence matrix $Y$ of interest point regions is taken as the objective function, and a matrix decomposition model is used to solve it. The new interest point score is represented by Eq. (3).

$$\hat{r}_{ul} = p_u \cdot q_l + x_u \cdot y_l \tag{3}$$

In Eq. (3), $q_l$ represents the matrix of interest points. $p_u$ is the implicit vector of tourists. In practical recommendations, tourists are easily influenced by social circle factors, and similar preferences between tourists and friends can easily lead to the final target selection. Therefore, it is necessary to calculate the similarity between them [15], which is represented by Eq. (4).

$$s(u,v) = \theta \times \hat{F}_{uv} + (1-\theta) \times \frac{|F_u \cap F_v|}{|F_u \cup F_v|} \tag{4}$$

In Eq. (4), $\theta$ is the adjustment parameter. $\hat{F}_{uv}$ means the friend relationship judgment. $F_u$ represents a collection of tourist friends. $F_v$ represents the collection of friends of user $v$. The objective loss function of the TSAS model is obtained by integrating the regional division, tourist social factors, tourist activity factors, and interest point influencing factors into the traditional matrix factorization model, which is represented by Eq. (5).

$$\lim_{P,Q,X} \sum_{(u,L) \in (U,L)} (r_{ul} - \hat{r}_{ul})^2 + \frac{\lambda_1}{2}(\|p_u\|^2 + \|q_l\|^2)$$
$$+ \frac{\lambda_2}{2}\|x_u\|^2 + \frac{\lambda_3}{2} \sum_{v \in F_u} s(u,v) \times \|p_u - p_v\|^2 \tag{5}$$

In Eq. (5), $\lambda_1$, $\lambda_2$, and $\lambda_3$ are the weights controlled by three factors. $p_v$ is the popularity of the user's area. $r_{ul}$ represents the actual point of interest score. $Q$ represents an implicit vector of interest points. $p$ represents the implicit vector of tourists. In order to improve the training effect of the objective function, the gradient descent method is used to optimize the parameters of the objective function. The gradient of $p_v$ is represented by Eq. (6).

$$\frac{\partial \Theta}{\partial p_v} = -2(r_{ul} - \hat{r}_{ul}) \cdot s(u,v) \cdot q_1 - \lambda_3 \cdot s(u,v) \cdot (p_u - p_v) \tag{6}$$

In Eq. (6), $\frac{\partial \Theta}{\partial p_v}$ represents the gradient of $p_v$. The gradient of $q_l$ is represented by Eq. (7).

$$\frac{\partial \Theta}{\partial q_l} = -2(r_{ul} - \hat{r}_{ul}) \cdot (p_u + \sum_{v \in F_u} s(u,v)p_v) + \lambda_2 q_1 \tag{7}$$

In Eq. (7), $\frac{\partial \Theta}{\partial q_l}$ represents the gradient of $q_l$. The gradient of $x_u$ is represented by Eq. (8).

$$\frac{\partial \Theta}{\partial q_l} = -2(r_{ul} - \hat{r}_{ul}) \cdot y_l + \lambda_2 x_u \tag{8}$$



Fig. 3. Recommended process for rural tourist attractions.

By using Eq. (6) to Eq. (8), the gradient of each parameter can be obtained. TSAS continuously optimizes the parameters in each iteration until the model converges or reaches the maximum number of iterations, completing the model training. Fig. 3 shows the entire recommendation model.

### B. Construction of a Recommendation Model Based on Tourist Check-in Area and Time Factors

In the construction of traditional interest point recommendation system models, it is impossible to avoid data sparsity and implicit feedback problems in the system. To avoid data sparsity, further mining can usually be done on regional geographic location, topic categories, time series, and tourist social information. However, there are hidden user behaviour patterns in tourist check-in data, and effectively extracting contextual information hidden in tourist check-in data is the key to improving model recommendation effectiveness [16]. Therefore, the TSAS recommendation model is improved by introducing a greedy clustering algorithm to search for tourist check-in center locations and divide them into different regions based on check-in points, analyzing the impact of different regions on tourist check-in interests. Meanwhile, the sequence of tourist interest points during a certain period is analyzed. By analyzing the time to reflect the transfer characteristics of tourist interest points during a certain period, the impact of time factor on tourist check-in can be obtained [17]. Greedy clustering method is used to partition and confirm the regions to find the center of each region in the sparse tourist check-in data. Fig. 1 shows its schematic diagram.

According to Fig. 4, the greedy clustering method is used to sort the check-in times of interest points. The region with the most check-in times is selected as the center, and the selected region center is scanned again, with the region less

than the distance d as the center point, and placed in the region set. If the current check-in reaches the set threshold ratio, the area will be divided, and the check-in times in the center of the high area will decrease towards the surrounding areas. The division of tourist check-in areas is made more reasonable by using the above methods [18]. The set of tourist check-in centers is defined as $C_u$, and the probability of tourists arriving at a given point of interest $l$ is expressed using a central Gaussian model, represented by Eq. (9).

$$P(l \mid C_u) = \sum_{c_u=1}^{C_u} \frac{1}{dst(l,c_u)} \cdot \frac{f_{c_u}}{\sum_{i \in C_u} f_i} \qquad (9)$$

In Eq. (9), $\dfrac{1}{dst(l,c_u)}$ represents the distance from the point of interest to the center of the region. $f_{c_u}$ represents the check-in frequency in different regional centers. The check-in of tourists is inversely proportional to their distance, with the closer they are, the more they check in. Therefore, the closer tourists are to the center of $l$, the higher the probability of check-in. In addition to analyzing the impact of check-in centers, it is also necessary to consider the influence of time factors on tourist interests. Therefore, the time proximity method is adopted to divide tourist check-in into implicit tourist vectors and implicit interest point vectors, and the product of the two factors is used to fit the rating prediction matrix [19]. If the probability of tourists checking in at $l$ is defined as $p(F_{ul})$, then Eq. (10) can be obtained.

$$P(F_{ul}) \propto P^T Q \qquad (10)$$



Fig. 4. Division of tourist check-in centers.

In Eq. (10), $P^T$ represents the implicit matrix of tourists. $\propto$ is the fitting symbol. The check-in points of tourists have sequential characteristics in time, indicating that the check-in is influenced by time factors. For example, the check-in area at noon is concentrated in the restaurant area, and the check-in area in the morning or afternoon is concentrated in specific scenic areas. The $k$ interest points of tourists who check-in within a certain period of time are analyzed to effectively analyze the check-in patterns of tourists during a certain period of time. These points are regarded as the current check-in records. $N_k^T(l)$ is recorded as time neighbors. The implicit vectors of time neighbors are accumulated as the implicit vectors of tourist check-in interest points, represented by Eq. (11).

$$l^T = \frac{1}{k} \cdot \sum_{l' \in N_k^T(l)} l_{l'} \qquad (11)$$

Based on the above research, the transfer pattern of tourist interests during a certain period of time in Eq. (12) can be obtained.

$$P(F_{ul}) \propto (p_i \cdot q_j + p_i \cdot \frac{1}{k} \sum_{l' \in N_k^T(l)} l_{l'}) \qquad (12)$$

In Eq. (12), $P(F_{ul})$ represents the probability of tourist check-in for fusion time transfer. $q_j$ represents the implicit interest point $j$. $p_i$ represents the implicit vector of user $i$. By using Eq. 12), the transfer pattern of tourists' interest at a certain time can be obtained. In order to better improve the recommendation effect of the model, a probability matrix decomposition model is used to obtain the objective optimization function, represented by Eq. (13).

$$\min_{U,L} \sum_{i=1}^{U} \sum_{j=1}^{L} I_{ij}(h(F_{ij}) - h(p_i^T q_j))^2 + \lambda_1 \|p_i\|_F^2 + \lambda_2 \|q_j\|_F^2 \qquad (13)$$

In Eq. (13), $U$ represents the implicit matrix of tourists. $L$ means the implicit matrix of interest points. $I_{ij}$ is the attendance record of user $i$ at the point of interest $j$. $h(.)$ refers to the logistic function. $F_{ij}$ represents the attendance status of user $i$ for interest point $j$. $F$ is a set of quantities. The time-transfer characteristics of tourists are integrated with the check-in modeling features to obtain the probability value of the final tourist $u$ at interest point $l$, represented by Eq. (14).

$$P_{ul} = P(F_{ul}) \cdot P(l \mid C_u) \qquad (14)$$

Recall (R) and Precision (P) are introduced to effectively evaluate the practical application effect of the recommendation model, represented by Eq. (15).

$$\begin{cases} recall@k = \frac{1}{|U|} \sum_{u=1}^{|U|} \frac{|S_u(k) \cap V_u|}{|V_u|} \\ presion@k = \frac{1}{|U|} \sum_{u=1}^{|U|} \frac{|S_u(k) \cap V_u|}{k} \end{cases} \qquad (15)$$



Fig. 5. Construction process of rural tourism attraction recommendation system.

In Eq. (15), $k$ represents the number of recommended points of interest. $S_u(k)$ indicates the interest points recommended by the top $k$ tourists. $V_u$ is a collection of interest points that tourists truly visit. High R and P indicate high accuracy of the recommendation. Fig. 5 shows the entire model construction process. In addition, the constructed recommendation system will fully take into account the distance and time factors of tourists, which provides more suitable rural tourism attractions for tourists. In a recommendation system, the main factor options to consider will be set, including time distance, cost, and comprehensive experience factors. Tourists prioritize experience and will overlook factors such as time, distance, and cost. Their overall planning focuses more on experience and functionality. In terms of time distance, more attention is paid to the travel time and distance, while also taking into account the experience [20]. Cost is mainly considered based on cost-effectiveness, taking into account factors such as distance and time experience, to meet the quality of tourist experience as much as possible while reducing the cost of the visitor. Fig. 5 shows the construction process of the entire rural tourism attraction recommendation system.

## IV. SIMULATION TESTING OF ALGORITHMS

This section consists of two parts: model performance and practical scenario application analysis. The performance analysis part mainly tests the performance of the model on a universal dataset. In the actual scenario, specific rural tourism data are selected for training to test the application effect of different models in rural tourism attraction recommendation.

### A. Performance Analysis of Rural Tourism Recommendation Models

Experimental tests were conducted on the WINDOWS 10 64 bit platform to test the performance of the proposed rural tourism recommendation model, with a running memory of 64GB, an Intel i9 16 core processor, and a graphics card NVIDIA RTX4080. Simulation experiments were performed on the Matlab platform for analysis. The Gowalla and Yelp datasets were selected for the experiment. Gowalla has 32510 points of interest and 18737 users. Yelp has 30887 users and 18995 points of interest. Singular Value Decomposition (SVD) and Probability Matrix Factorization (PMF) models were introduced as recommended testing benchmarks. In actual testing, $\lambda_1$, $\lambda_2$, and $\lambda_3$ are important weight parameters that affect the training of the proposed model. Therefore, it is necessary to select appropriate regularization parameters for testing. Root Mean Squared Error (RMSE) was used to reflect the results in Fig. 6.

In Fig. 6, $\lambda_1$ is mainly responsible for weighting the implicit vectors of tourists and interest points. In Gowalla, when $\lambda_1$ was 0.5, RMSE was the lowest. In Yelp, when $\lambda_1$ was 0.3, RMSE was the lowest. Overall, the analysis shows that Yelp is relatively sparse, and the model performs best when the dataset is sparse with a $\lambda_1$ of 0.3. $\lambda_2$ mainly affects the weights in the tourist activity matrix. Through experimental analysis, in Gowalla, the best model performance was achieved when $\lambda_2$ was 0.3. When the Yelp was sparse, the best model performance was achieved when $\lambda_2$ was 5. $\lambda_3$ is a parameter that controls the social weight of tourists. In Gowalla and Yelp, the best performance was achieved when $\lambda_3$ was 0.3 and 0.6, respectively. Therefore, in subsequent experiments, effective weights are set based on the sparsity of the test samples to ensure the testing performance of the model. Meanwhile, $\theta$ represents the similarity adjustment parameter, which also has a direct impact on model testing in Fig. 7.

In Fig. 7, regardless of whether the dataset was sparse or not, $\theta$ had no significant impact on the performance of the model. When $\theta$ was 0.5, the model had the best testing performance. Therefore, based on the above experimental results, appropriate parameter values were selected for comparison. Fig. 8 shows the comparison results of recommendation accuracy between different models.

According to the results in Fig. 8, in Gowalla, the proposed model achieved the earliest convergence and had the highest recommendation accuracy of 0.965, while the PMF and SVD recommendation models were 0.912 and 0.946, respectively. Meanwhile, in Yelp, the recommended accuracy of the proposed model, SVD, and PMF was 95.65, 0.832, and 0.795, respectively. When the dataset is sparse, the recommendation performance of PMF and SVD significantly decreases, while the proposed model still has excellent recommendation performance. Fig. 9 compares the errors of two models.

According to Fig. 9, in Gowalla, the RMSE loss of PMF, SVD, and the proposed model towards convergence was 0.425, 0.335, and 0.120, respectively. In Yelp, when PMF, SVD, and the proposed model tended to converge, the RMSE loss was 0.865, 0.432, and 0.132, respectively. The proposed model has lower overall RMSE loss and better performance.

Fig. 6. Comparison of different weight parameters.



Fig. 7. Comparison of different similarity adjustment parameters.

Fig. 8.    Comparison of recommendation accuracy among different recommendation models.



Fig. 9.    Comparison of error performance of different recommendation models.

### B. Analysis of Practical Application Scenarios of Tourism Recommendation Models

Crawler technology was used to crawl Ctrip tourist comment information, including 215654 rural tourism check-in score data, catering data, etc. Baidu Map platform was used to search for the longitude and latitude coordinates of rural tourist attractions, and the 8km range of scenic spots were classified into the same section. Finally, 289456 distance section data were collected, and the final regional feature data of rural scenic spots were obtained through sorting. Table I shows the specific parameters.

TABLE I.    COLLECTION OF INFORMATION FOR RURAL TOURISM DESTINATIONS

| Types of | Training set | Test set | Section data |
|---|---|---|---|
| File size | 2.26M | 2.41M | 5.07M |
| Number of tourists | 23156 | 29651 | 32546 |
| Number of attractions | 16581 | 19635 | 19465 |
| Number of elements | 144488 | 71165 | 289456 |
| Minimum score/minimum height range | 1 | 1 | 0 |
| Maximum score/maximum pitch range | 5 | 5 | 425 |
| The proportion of the dataset | 67.00% | 33.00% | - |

In Fig. 10, SVD and SVD models are still used as benchmark models, and the recommendation performance of the models in actual rural scenic spots is compared.

According to Fig. 10, when the number of recommendations was 5, all three models had the best recommendation performance. The accuracy values of PMF, SVD, and the proposed model were 0.098, 0.111, and 0.138, respectively, when the recommended quantity was 5. Simultaneously comparing the recall rates of different models, when the number of recommendations was 5, the recall rate of PMF, SVD, and the proposed model was 0.048, 0.051, and 0.069, respectively. The proposed model has better accuracy and recall performance than other recommendation models. Finally, Fig. 11 compares the itinerary planning effects of three models in rural tourism scenarios.

Fig. 11 shows the effect of travel itinerary planning for different models, which include three recommendation modes: time distance, cost, and comprehensive experience. Travel arrangements are planned according to the needs of the tourists. As tourists spent more time in rural scenic areas, the planning time of different models varied significantly. Among them, the overall planning time of the PMF model was the longest, with the highest planning time reaching 11200ms after the tourist travel time reached 330 minutes. The longest planning time for SVD was 6212ms. The best performing model is the proposed one. Although the planning time of the proposed model increased after the tourist's play time reached 330 minutes, the planning efficiency was still the highest compared with the other two models, and the longest planning time of the model was 1956 ms. Therefore, the proposed model has excellent rural tourism recommendation performance.

Fig. 10. Comparison of accuracy and recall performance of different models.



Fig. 11. Comparison of tourism itinerary planning efficiency among different models.

## V. CONCLUSION

Rural tourism has attracted a large number of tourists due to its unique culture and characteristics, but the recommendation of rural scenic spots has always faced difficulties and couldn't meet practical needs. A region segmentation-based recommendation technique is proposed for this purpose. Firstly, tourist check-in and geographical location are considered, and the check-in situation is used to reflect the process of segmenting regions, thereby achieving content recommendation. While in practical content recommendation, data sparsity and visitor interest transfer issues need to be considered as well. Therefore, the modeling is based on tourist check-in areas and time factors to capture the temporal changes of tourists. Finally, different itinerary planning schemes are matched according to the needs of tourists to achieve recommendations and itinerary planning for rural scenic spots. In the experimental analysis of model performance, the proposed model, SVD, and PMF models achieve recommendation accuracy of 0.956, 0.832, and 0.795 in Yelp, respectively. Meanwhile, the RMSE loss is 0.865, 0.432, and 0.132, respectively, when PMF, SVD and the proposed model tend to converge. In practical scenario application analysis, the optimal recall rate of the proposed model is 0.069, and the PMF and SVD are 0.048 and 0.051, respectively. Comparing the travel planning efficiency of different models, the highest time consumption of the proposed model is 1956ms, while PMF and SVD are 11200ms and 6212ms, respectively. Therefore, the proposed model has excellent recommendation and itinerary planning effects in rural tourism attraction recommendation. There are still

shortcomings in this study. The proposed method relies on tourist check-in data for recommendation. Although it introduces time-sensitive nearest weight, in some cases, the recommendation system may still not be able to fully capture the instantaneous interest changes of users. In the future, research technology also needs to consider regional meteorological factors, holidays and other factors, fully considering the impact of these factors on tourist recommendations, to optimize the practical application effect of recommendation technology.

### REFERENCES

[1] Wei G, Wu Q, Zhou M. A hybrid probabilistic multiobjective evolutionary algorithm for commercial recommendation systems. IEEE Transactions on Computational Social Systems, 2021, 8(3): 589-598.

[2] Cao B, Zhang Y, Zhao J, Liu X. Recommendation based on large-scale many-objective optimization for the intelligent internet of things system. IEEE Internet of Things Journal, 2021, 9(16): 15030-15038.

[3] Cao B, Zhao J, Lv Z, Yang P.Diversified personalized recommendation optimization based on mobile data. IEEE Transactions on Intelligent Transportation Systems, 2020, 22(4): 2133-2139.

[4] Cui Z, Xu X, Fei X U E, Cai X, Cao Y. Personalized recommendation system based on collaborative filtering for IoT scenarios. IEEE Transactions on Services Computing, 2020, 13(4): 685-695.

[5] Zhou X, Li Y, Liang W. CNN-RNN based intelligent recommendation for online medical pre-diagnosis support. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2020, 18(3): 912-921.

[6] Cho J, Ahmed S, Hilbert M, et al. Do search algorithms endanger democracy? An experimental investigation of algorithm effects on political polarization. Journal of Broadcasting & Electronic Media, 2020, 64(2): 150-172.

[7] Nitu P, Coelho J, Madiraju P. Improvising personalized travel recommendation system with recency effects. Big Data Mining and Analytics, 2021, 4(3): 139-154.

[8] Giglio S, Bertacchini F, Bilotta E, et al. Machine learning and points of interest: typical tourist Italian cities. Current Issues in Tourism, 2020, 23(13): 1646-1658.

[9] Huang F, Xu J, Weng J. Multi-task travel route planning with a flexible deep learning framework. IEEE Transactions on Intelligent Transportation Systems, 2020, 22(7): 3907-3918.

[10] Wang S, Gong M, Wu Y, Zhang M. Multi-objective optimization for location-based and preferences-aware recommendation. Information Sciences, 2020, 513(6): 614-626.

[11] Al Fararni K, Nafis F, Aghoutane B, et al. Hybrid recommender system for tourism based on big data and AI: A conceptual framework. Big Data Mining and Analytics, 2021, 4(1): 47-55.

[12] Asselman A, Khaldi M, Aammou S. Enhancing the prediction of student performance based on the machine learning XGBoost algorithm. Interactive Learning Environments, 2023, 31(6): 3360-3379.

[13] Fan W, Ma Y, Li Q, Cai G. A graph neural network framework for social recommendations. IEEE Transactions on Knowledge and Data Engineering, 2020, 34(5): 2033-2047.

[14] Dhelim S, Ning H, Aung N, Huang R. Personality-aware product recommendation system based on user interests mining and metapath discovery. IEEE Transactions on Computational Social Systems, 2020, 8(1): 86-98.

[15] Huang Z, Xu X, Zhu H, Zhou MC. An efficient group recommendation model with multiattention-based neural networks. IEEE Transactions on Neural Networks and Learning Systems, 2020, 31(11): 4461-4474.

[16] Grossman G, Kim S, Rexer J M. Political partisanship influences behavioral responses to governors' recommendations for COVID-19 prevention in the United States. Proceedings of the National Academy of Sciences, 2020, 117(39): 24144-24153.

[17] Liang S. Research on Route Planning of Red Tourist Attractions in Guangzhou Based on Ant Colony Algorithm. Automation and Machine Learning, 2023, 4(1): 8-16.

[18] Zhou X, Li Y, Liang W. CNN-RNN based intelligent recommendation for online medical pre-diagnosis support. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2020, 18(3): 912-921.

[19] Liu K, Sun Y, Yang D. The Administrative Center or Economic Center: Which Dominates the Regional Green Development Pattern? A Case Study of Shandong Peninsula Urban Agglomeration, China. Green and Low-Carbon Economy, 2023, 1(3), 110–120.

[20] Ly A, El-Sayegh Z. Tire wear and pollutants: An overview of research. Archives of Advanced Engineering Science, 2023, 1(1): 2-10.

# A Hybrid GAN-BiGRU Model Enhanced by African Buffalo Optimization for Diabetic Retinopathy Detection

Dr. Sasikala P[1], Sushil Dohare[2], Dr. Mohammed Saleh Al Ansari[3], Janjhyam Venkata Naga Ramesh[4],
Prof. Ts. Dr. Yousef A.Baker El-Ebiary[5], Dr. E.Thenmozhi[6]

Assistant professor, Department of Computer Science,
Government Science college (Nrupathunga University), Bangalore, Karnataka[1]
Department of Epidemiology, College of Public Health and Tropical Medicine, Jazan University, Saudi Arabia[2]
Associate Professor, College of Engineering, Department of Chemical Engineering, University of Bahrain, Bahrain[3]
Assistant Professor, Department of Computer Science and Engineering,
Koneru Lakshmaiah Education Foundation, Vaddeswaram, Guntur Dist., Andhra Pradesh - 522302, India[4]
Faculty of Informatics and Computing, UniSZA University, Malaysia[5]
Associate Professor, Department of Information Technology, Panimalar Engineering College, Chennai, India[6]

*Abstract*—Diabetic retinopathy (DR) is a severe complication of diabetes mellitus, leading to vision impairment or even blindness if not diagnosed and treated early. A manual inspection of the patient's retina is the conventional way for diagnosing diabetic retinopathy. This study offers a novel method for the identification of diabetic retinopathy in medical diagnosis. Using a hybrid Generative Adversarial Network (GAN) and Bidirectional Gated Recurrent Unit (BiGRU) model, further refined using the African Buffalo Optimization algorithm, the model's capacity to identify minute patterns suggestive of diabetic retinopathy is improved by the GAN's skill in extracting complex characteristics from retinal pictures. The technique of feature extraction plays a critical role in revealing information that may be hidden yet is essential for a precise diagnosis. Then, the BiGRU part works on the characteristics that have been extracted, efficiently maintaining temporal relationships, and enabling thorough information absorption. The combination of GAN's feature extraction capabilities with BiGRU's sequential information processing capability creates a synergistic interaction that gives the model a comprehensive grasp of retinal pictures. Moreover, the African Buffalo Optimization technique is utilized to optimize the model's performance for improved accuracy in the identification of diabetic retinopathy by fine-tuning its parameters. The current study, which uses Python, obtains a 98.5% accuracy rate and demonstrates its amazing ability to reach high levels of accuracy in Diabetic Retinopathy Detection.

*Keywords*—*African Buffalo Optimization (ABO); Bidirectional Gated Recurrent Unit (BI-GRU); Generative Adversarial Network (GAN); diabetic retinopathy; medical diagnosis*

## I. INTRODUCTION

Diabetic retinopathy is a significant and escalating public health concern that affects individuals with diabetes, representing a leading cause of blindness worldwide [1] .This progressive eye disease is primarily attributed to prolonged exposure to elevated blood sugar which can lead to impairment of the blood vessels in the retina, the light-sensitive tissue at the back of the eye. These blood vessels may start to leak blood or other fluids into the eye as they degrade, which can cause retinal edema, blurred vision, and eventually vision loss [2]. Regular screening and early identification are essential for successful intervention and management of diabetic retinopathy since the severity of the condition might vary, and its development may be asymptomatic in its early stages. Overall, ophthalmologists' dilated eye exams have been the gold standard for detecting diabetic retinopathy [3]. While these tests are still necessary, the rising incidence of diabetes and the danger of diabetic retinopathy that goes along with it highlight the need for effective and scalable diagnostic measures [4]. Medical imaging and artificial intelligence (AI) have made amazing strides in recent years, and these developments are increasingly being used to transform the identification of diabetic retinopathy [5]. There are many phases of diabetic retinopathy, from non-proliferative retinopathy, which is milder, to proliferative retinopathy, which is more severe and involves the development of aberrant blood vessels on the retina. Early identification is important because it allows for appropriate management to stop or postpone vision loss, preserving patients' visual function and general quality of life. Diagnosing diabetic retinopathy has mostly included dilated eye exams performed by qualified ophthalmologists. In these tests, dilating eye drops are used to provide a thorough view of the retina, which is followed by a careful visual assessment [6]. Despite being extremely precise, this procedure is time- and resource-consuming, and human interpretation might vary. It also presents a major accessibility and scalability barrier, particularly in areas with poor access to eye care professionals. The development of non-invasive, high-resolution imaging methods for the retina is the result of developments in medical imaging [7] They provide precise observation of retinal structures and abnormalities, such as Ophthalmic coherence tomography, and fundus photography (OCT), and fluorescein angiograph. Although these methods of imaging have increased diagnostic precision and allowed

for early diagnosis, they continue to depend on human judgment and are prone to differences in skill

Huge collections of retinal pictures are being used to models so they can identify the telltale symptoms of diabetic retinopathy [8]. These DL-driven solutions have shown outstanding diagnostic speed, accuracy, and consistency; they also possess the potential to supplement current approaches and reduce access to care gaps [9]. The conditions of detecting diabetic retinopathy may change as a result of the combination of DL and medical imaging [10]. Automation of image processing has the potential of expanding screening accessibility, easing the strain on healthcare resources, and enhancing diagnosis precision, particularly in areas where access to ophthalmologists is constrained [11]. The ramifications of these developments, particularly their effect on patient outcomes and healthcare delivery, will be explored in this study. With an emphasis on how this convergence might help preserve sight and improve the quality of life for people with diabetes, researchers will also examine existing research, difficulties, and the future potential of DL in diabetic retinopathy diagnosis. The main issue with diabetic retinopathy is that it progresses initially without any symptoms. Early identification is crucial because patients may not become aware of the disease until it has progressed to an advanced stage. From moderate non-proliferative retinopathy to severe proliferative retinopathy, the disorder progresses through numerous phases and is defined by the growth of aberrant blood vessels in the retina [12]. Effective management can stop or delay vision loss, preserving the quality of life for those with diabetes, therefore early detection of diabetic retinopathy is essential (Lin et al. 2021). While conventional diagnostic methods are dependable, they can also be resource-intensive and not always accessible, particularly in places where there are few or no access to eye care practitioners. As a result, investigating cutting-edge methods like artificial intelligence and enhanced medical imaging has become a potential direction to solve this important healthcare issue [13].Technology is becoming increasingly important in the transformation of diabetic retinopathy diagnosis. [14]. In order to provide facts about the history, current, and destiny of this crucial field of healthcare, this investigation aims to shed light on the changing landscape of diabetic retinopathy detection.An innovative method that combines the benefits of Generative Adversarial Network (GAN) with Bidirectional Gated Recurrent Units (BiGRUs) for improved image processing is called the African Buffalo Optimization Guided Hybrid GAN-BiGRU. GAN are excellent at handling spatial data, but BiGRUs are good at handling sequential data. This study presents a possible structure for a wide range of applications, promising enhanced accuracy and efficiency in image analysis and pattern identification by combining these two deep learning paradigms and the African Buffalo Optimization method.

The proposed research aims to address these limitations by introducing "A Hybrid GAN-BiGRU Model Enhanced by African Buffalo Optimization for Diabetic Retinopathy Detection." The significance of this research lies in its potential to provide a more efficient, accessible, and accurate method for the early detection and diagnosis of diabetic

retinopathy. The gap in existing knowledge that the study aims to fill is the development of a novel approach that combines advanced neural network architectures (GAN and BiGRU) with a bio-inspired optimization technique (African Buffalo Optimization) to enhance the precision and efficiency of diabetic retinopathy detection.

The following are the study's key contributions.

- ➢ Proposes a hybrid Generative Adversarial Network (GAN) and Bi-directional Gated Recurrent Unit (BiGRU) model for diabetic retinopathy detection.

- ➢ Implements the African Buffalo Optimization algorithm to refine and optimize the model's performance, enhancing its efficiency.

- ➢ Enhances the identification of subtle patterns associated with diabetic retinopathy by leveraging GAN's effective feature extraction capabilities from retinal images.

- ➢ Highlights the crucial role of feature extraction in revealing hidden information essential for precise diagnoses in diabetic retinopathy.

- ➢ Demonstrates the contribution of Bi-directional Gated Recurrent Unit (BiGRU) in maintaining temporal relationships and facilitating thorough information absorption from the extracted characteristics.

- ➢ Explores the synergistic interaction between GAN's feature extraction and BiGRU's sequential

Overall, the study contributes to advancing the field of diabetic retinopathy diagnosis by combining advanced neural network architectures and bio-inspired optimization techniques for precise and efficient detection. The format for the enduring paragraphs is as follows: The relevant work based on various methodologies for diabetes prediction is examined in Section II, and the research gap is identified in section III. The feature selection and classification process for the proposed method is explained in the Section IV. The outcomes and considerations are covered in Section V; the prospective applications for the future are covered in Section VI.

## II.    RELATED WORKS

The disease known as diabetic retinopathy (DR), which obliterates the retinal veins, can cause blind. To diagnose this lethal illness, colored fundus injections are frequently used. The manual examination of the aforementioned photos (by physicians) is tedious and prone to mistakes. As a result, a variety of computer vision engineering approaches are used to forecast the DR's appearances and phases autonomously.

These techniques can't properly categorize DR's various phases since they are operationally costly and don't retrieve extremely nonlinear information. In order to hasten the training of models and convergence, Khan et al. [15] focuses on categorizing the DR's several phases by means of the minimum constraints that may be learned feasibleThe VGG-NiN paradigm is built by stacking the VGG16, the spatial pyramid layer for pooling (SPP), and the network-in-network (NiN). It is a highly nonlinear scale-invariant deeper method. The recommended VGG-NiN device is capable of handling a DR image of any dimension in along with the benefits of the SPP layer. The framework also gains extra nonlinearity from the stacking of NiN, which enhances classification overall. The recommended strategy beats cutting-edge approaches when it comes to of efficiency and effective use of computer resources, according to data from experiments. The model's construction and preprocessing techniques need to be changed to boost output.

Identifying diabetic retinal disease in its early stages and predicting the potential presence of Micro aneurysms in fundus images were incredibly challenging for a long time. Long-term high blood glucose levels cause diabetic retinopathy (DR), which. The field of deep learning is quickly advancing, results in micro vascular problems and permanent blindness. The initial indications of DR are the development of micro aneurysms and macular edema in the retina, and prompt detection can decrease the chance of developing non-proliferated diabetic retinopathymaking it an effective method for offering an intriguing answer to difficulties involving clinical picture interpretation. (Qiao, Zhu, and Zhou [16] proposed system analyzes the existence of a micro aneurysm in a fundus image via convolutional neural network techniques that incorporate a deep learning approach as an essential element and are accelerated by a GPU (Graphics Processing Unit). This will enable outstanding performance and low-latency inference for the identification and segmentation of medical images. The fundus image is categorized as normal or diseased using the semantic segmentation technique The process of semantic segmentation divides the image pixel into groups based on their common semantics in order to identify a micro aneurysm's characteristics. This gives ophthalmologists a computerized method to help them classify fundus pictures as quickly, mild, or extreme NPDR. The early identification and prognosis method for non-proliferative diabetic retinopathy was suggested, and it has the ability of successfully developing a deep convolution neural network (CNN) for semantic division of retina images, which can improve the efficacy and precision of NPDR (non-proliferated diabetic retinopathy) prediction. It is essential to conduct testing on various datasets and in real-life healthcare environments to determine the system's usefulness and dependability.

Among the conditions that poses the greatest risk to vision is retinopathy caused by diabetes (DR), a consequence brought on by elevated blood sugar levels. However, an ophthalmologist must manually collect DR screening, which is time-consuming and subject to error. The enormous rise in the number of diabetic patients has led to an emphasis on automated DR diagnosis in recent years. Additionally,

Convolutional Neural Networks (CNN) have proven themselves to be state-of-the-art for DR stage diagnosis in recent times. (Farag, Fouad, and Abdel-Hamid [17]offers a fresh, a system that autonomously learns how to calculate brightness from just one Colour Fundus Photograph (CFP). The suggested method builds a visual embedding using DenseNet169's encoder. Convolutional Block Attention Module (CBAM) is also added on the highest of the encoder to boost its ability to discriminate. The algorithm is then trained using the Kaggle Asia Pacific Tele-Ophthalmology Society (APTOS) database utilizing cross-entropy loss. .THIS approach makes a substantial contribution by accurately classifying the degree of diabetic retinopathy intensity while requiring less time and spatial difficulty, making it a potential contender for autonomous diagnosis. the effectiveness of various CBAM designs. It is recommended to apply several unbalanced learning algorithms, and expanding the dataset will improve results.

Hemanth, Deperlioglu, and Kose [18] provide an innovative hybrid strategy for the retinal fundus imaging-based diagnosis of retinopathy caused by diabetes. To improve diagnosis accuracy, our mixed strategy specifically integrates processing of images and deep learning. It is well known that manually interpreting these photos is a laborious, time-consuming operation requiring substantial knowledge. Medical experts turn to computer vision systems for assistance in tackling this problem, and intelligence diagnostic techniques have grown in importance. Using image processing methods like equalization of histograms and contrast-limited adaptive equalization of the histogram in the present investigation, we suggest a diagnostic method using convolutional neural networks. It is essential to enhance integrating additional imaging modalities, conducting large-scale clinical validation, and developing real-time monitoring capabilities.

To determine whether certain actions may be made to improve efficiency and solution quality, use a variety of image processing techniques. A prompt evaluation and therapy are required for diabetic retinopathy in order to prevent visual loss. Since they are concealed in minute and subtle shapes beneath the eye's structure, the medical condition's lesions are challenging to detect. Maaliw et al. [19] built an efficient process to extract pertinent features utilizing a variety of preparation methods, a visual segmentation design (DR-UNet) that has an impressive spatial pyramid pool, and an attention-aware neural convolutional networks with multiple ResidualNet-based sections. Experimental findings show that our system's precision for segmentation is 87.10% (intersection over union) and 84.50 % (dice similarity coefficient). The gradual converging of training/validation loss and accuracy further supported the claim that the 99.20% classification performance outperformed prior systems. In order to more accurately diagnose the illness in both its early and severe phases, this investigation has the potential to complement conventional diagnostics. Determine the severity of DR and create an improved framework going forward.

Diabetics must find Diabetic Retinopathy (DR) early to reduce their chance of going blind. Numerous research show that Deep Convolutional Neural Network (CNN)-based

methods are efficient for enabling automated DR identification via categorization of patient retinal pictures. In order to assist their CNN training, these techniques often rely on an enormous data set made up of retinal pictures with predetermined categorization labels. Finding sufficient accurately labeled pictures to serve as model training examples, nevertheless, is not always simple. In addition, as a CNN gets deeper, training it takes greater time and is more probable to result in over fitting, particularly if you utilize a big training dataset. In order to categorize retinal pictures, it is important to investigate a more straightforward CNN-based method that is still successful on tiny data sets. W. Chen et al [20] proposes a method for classifying retinal images that integrates multi-scale shallow CNNs. Research on open-source datasets demonstrates that, when compared with existing representational combined CNN learning algorithms, the suggested method can increase classification accuracy by 3% for short datasets. In comparison to other typical techniques like conventional CNN, LCNN, and VGG16noFC on a larger dataset, the performance of the integrated shallow CNN model will be enhanced by the modification of picture samples and repeated dataset sampling.

The limitations faced by current methods in detecting and diagnosing diabetic retinopathy (DR). The existing approaches, including deep learning techniques like Convolutional Neural Networks (CNNs), encounter challenges in categorizing the various phases of DR due to operational costs and difficulty in capturing highly nonlinear information. These methods also vary in terms of efficiency, resource utilization, and accessibility, with some relying on specialized hardware and complex models. The need for extensive datasets, lack of model explainability, and the importance of clinical validation further hinder the development of accessible and accurate diagnostic tools, particularly in resource-constrained healthcare settings.

## III. PROBLEM STATEMENT

Diabetic retinopathy (DR) poses a significant risk to vision, leading to blindness if not diagnosed and treated promptly. The conventional manual examination of retinal images for DR diagnosis is time-consuming, error-prone, and relies heavily on the expertise of medical professionals [20]. Existing computer vision engineering approaches, including those utilizing deep learning techniques like Convolutional Neural Networks (CNNs), face challenges in properly categorizing the various phases of DR due to operational costs and difficulty in capturing highly nonlinear information [16]. Furthermore, the current methods vary in terms of efficiency, resource utilization, and accessibility [17], with some relying on specialized hardware and complex models. The necessity for extensive datasets, lack of model explainability, and the importance of clinical validation further hinder the development of accessible and accurate diagnostic tools, especially in resource-constrained healthcare settings. As the number of diabetic patients increases, there is a pressing need to overcome these limitations and devise more efficient, accessible, and accurate methods for the early detection and diagnosis of diabetic retinopathy, ultimately improving patient care and outcomes we proposes

## IV. PROPOSED AFRICAN BUFFALO OPTIMIZATION BASED HYBRID GAN -BiGRU

The proposed African Buffalo Optimization (ABO) based hybrid Generative Adversarial Network (GAN) and Bi-directional Gated Recurrent Unit (BiGRU) model represents a novel approach for diabetic retinopathy detection. ABO is integrated into the model to optimize its performance by fine-tuning parameters, enhancing its accuracy in identifying diabetic retinopathy. The hybrid architecture combines GAN's proficiency in extracting complex features from retinal images with BiGRU's sequential information processing capabilities. This synergistic interaction, augmented by ABO, result in a comprehensive and efficient model, showcasing the potential of bio-inspired optimization techniques in advancing the accuracy and reliability of diabetic retinopathy diagnosis. Fig. 1 demonstrates proposed method.



Fig. 1. Proposed ABO with GAN-BIGRU.

## A. Datasets

Data set received from the Kaggle coding website (https://www.kaggle.com), which comprises over 80,000 images, each with an approximate resolution of 6 million pixels and various scales of retinopathy. To effectively process this dataset images are resized and conducted our deep learning experiments using a high-end GPU, specifically the NVIDIA K40c. This GPU is equipped with 2,880 CUDA cores and is supported by the NVIDIA CUDA Deep Neural Network library (cuDNN) for GPU-accelerated machine knowledge. To facilitate the training process, This research leveraged the capabilities of the Keras deep learning package (http://keras.io/) with Theano as the machine learning backend (http://deeplearning.net/software/theano/). This choice was made due to the availability of excellent documentation and the advantage of shorter calculation times. The performance achieved was remarkable. In fact, our system can classify an image in a mere 0.04 seconds, providing the potential for real-time feedback to benefit patients [21].

## B. Image Processing

Image preprocessing in the context of image analysis involves a series of techniques to enhance image quality and make it more suitable for subsequent analysis. Histogram equalization is a specific method used to improve contrast in images. It works by redistributing pixel intensities across the entire range, resulting in an image with enhanced contrast, which can be particularly useful for improving the visual quality of images and making them more amenable to various image analysis tasks.

*1) Histogram equalization:* It remains one of among the most popular methods for enhancing the clarity and quality of photographs that go through processing. The goal image's histogram's range of motion is increased to accomplish this. The HE quickly converts the original image's irregular gray levels into the resulting image's consistent level of gray. The produced image has a homogeneous pattern of gray levels as a result. It is therefore reasonable to say that the HE is used to generate an even histogram. Finally, HE delivers a new intensities value for each pixel depending on its prior intensity level. The visual appeal of the picture is enhanced and its histogram is dispersed across a larger range since the histogram for low-contrast images is narrow and focused near the middle of the gray scale. The HE improves the brightness of the image by flattening and stretching the given data's histogram's range of motion.

The following theoretical approach to HE can help us understand greater detail regarding it:

It is feasible to think about a digital picture $I(i,j)$ with X total pixels and a grayscale between [0, K'-1]. Following there, an equation may be used to determine the density function probability of the associated picture in Eq. (1)

$$p(k') = \frac{n_{k'}}{X}, \text{ for k'=0,1,......K-1} \qquad (1)$$

where $n_{k'}$ is the overall amount of grayscale pixels  k n the picture. Additionally, Eq. (2) may be used to determine the image's $(i, j)$  cumulative distribution function as shown in Eq. (2)

$$B(k') = \sum_{m=0}^{k'} P'_m \text{ for k'=0,1,.....K-1} \qquad (2)$$

By taking into account the results of the cumulative distributions function. A level of input k is matched with a level of output k by $HE_{k'}$ may be used to do this in Eq. (3)

$$HE_{k'} = (K-1).B_{k'} \qquad (3)$$

As a consequence, eqn may be used to determine the gain $HE_{k'}$ at the output level for the typical HE in Eq. (4)

$$\Delta HE_{k'} = HE_{k'} - (HE_{k'}-1) = (K-1).P'(k) \qquad (4)$$

It is feasible to show that the rise in the level of  of $HE_{k'}$ is proportionate to the likelihood of the corresponding value k in the setting of the initial picture by taking into account the linked equations. Explains the HE procedure as it appears over the histogram data in short. For photos with widely spaced tonal zones, HE is highly helpful in enabling the observation of images with a very light backdrop and a dark foreground. Expanding out the disparity between neighboring regions, that enables making noticeable the variations within the processing regions, enables the HE to reveal concealed data inside a picture.

## C. Segmentation using Watershed Algorithm

The watershed method is guided by the probability maps, which improves the segmentation process' accuracy and capacity to discern among various sections of the picture. This method helps in the precise delineation of tumor borders in healthcare imaging applications and efficiently handles the problem of over-segmentation, which has become a prevalent issue in segmenting image assignments is given in Eq. (5) (6) and (7)

$$\text{Entity Part} = \{\text{Locality}|\ \text{Strength(Locality ) }\geq \text{threshold}\} \qquad (5)$$

$$\text{Threshold} = \min(\text{Optimistic}) + \big[\max(\text{Optimistic}) - \min(\text{Optimistic})\big] \times \text{Rate} \qquad (6)$$

$$\text{Optimistic} = \{\text{Intensity} \mid \text{Intensity} \geq \text{mean(Intensities)}\} \qquad (7)$$

'"Rate" is a number between 0 and 1, whereas "Optimistic" refers to the bright portion of an AO-SLO picture. Additionally, the backdrop markers were set to match the findings from local morphology processing, while the remaining region was assigned to identifiers for an unknown location. Utilizing the previously described contour-length threshold-based technique, the marker-controlled watershed segment technique was performed repeatedly at an accelerating pace till no conjunction-containing areas were found. With such cycles, every region's contour-length is below the cut-off value of one cone photoreceptor cell's contour-length, eliminating the conjunction-containing areas. Ultimately, all marker-controlled watershed segmentation rectangles are added to the outcomes of the regional morphology process.  provides an illustration of the outcomes of the watershed method applied to a typical picture patch; as

can be seen, the conjunction-containing patches are distinct [22].

*1) GAN-based retinal vasculature extraction:* A novel Pix2Pix Generative Adversarial Network (GAN) design was used in this investigation. This design, which was first presented by Ian J. Goodfellow in 2014, consists of the Generators and Discrimination sub-models. The Generator is the component that creates data samples while the Discriminator tries to distinguish between produced and actual data, putting both of these models in competition with one another. Training doesn't stop until the discriminating agent can't tell the two apart. In order to create a picture, the Generator network is given a fixed-length randomized seed noise, also known as a receptive vector. The generating method is built on top of this latent vector. For discrimination, the resultant picture and actual images are given into the discriminator. Following training, latent variables—which resemble locations within the issue's domain but are unable to be directly observed—are represented as multifaceted vector spaces termed latent spaces. High-level concepts from the raw data are captured in the latent space, which the algorithm uses to analyze events and generate fresh results.

As a model for categorization, the Discriminant separates created examples from actual samples based on the training data. In order to minimize the Discriminator's loss, the Generator and Discriminator losses are tracked throughout training. Fig. 2 shows the systematic architecture of GAN. Training improves the Discriminator's ability to discriminate between genuine and false data and the Generator's ability to produce accurate information. Upon reaching integration, the Generator automatically generates data that is almost realistic, and the Discriminator produces ½ for every input, making it unnecessary after training.

Applications for GANs may be found in many different fields, including pattern transfer, processing images, tracking traffic, and the creation of 3D objects. Visual translating, which converts a given input image into an output picture, is one important use.

There are several varieties of GANs, such as DCGAN, cGAN, Cycle GAN, and Info GAN. For picture further sampling, DCGAN employs transposition convolutional neural networks and deep convolutional nets. cGANs are appropriate for translating images to images because they let the GAN be conditionally trained using labels for classes. Similar assignments can be completed by Cycle GAN, however it can also use mismatched dataset for learning visual mappings. Comprehensible and significant representations can be learned using info GANs. A Pix2Pix GAN, a particular instance of cGAN that is frequently employed for translating images into images research, was employed in this investigation [23].

A cGAN is able to comprehend how to map from a perceived picture (x) and a random noise vector (z) to a generated image (y), expressed as Eq. (8)

$$G' = X', Z' = Y' \qquad (8)$$

The cGAN's loss coefficient is shown below Eq. (9)

$$\mathcal{L}'_c \text{GAN'(G',D')} = E'_{(X',Y')}[\log D'(X',Y')] + E'_{(X',Z')}[\log(1 - D'\big(X', G'(X',Z')\big))] \qquad (9)$$

In this case, the discrimination coefficient D tends to raise the previously indicated function, whereas the generator G tends to diminish it. An unconditionally version is applied to the GAN loss in order to compute the consequences of conditioning D, as can be shown below in Eq. (10)

$$\mathcal{L}\text{GAN'(G',D')} = E'_{(Y')}[\log D'(Y')] + E'_{(X',Z')}[\log(1 - D'\big(G'(X',Z')\big))] \qquad (10)$$

Throughout the further sampling and downsampling processes, the Generator in the Pix2Pix GAN forms a UNet structure using a Resnet. To further reduce distortion, a reduction function is implemented in G as follows $\mathcal{L}'_{L1'}$ in Eq. (11)

$$\mathcal{L}'_{L1'(G')} = E'_{(X',Y',Z')}[[\|\text{Y'-G'}(X',Z')\|]] \qquad (11)$$

Having a patch size of 70 x 70, the Discriminators is a patchGAN. The Pix2Pix GAN's ultimate loss value is represented by a calculation that combines the cGAN loss with the $\mathcal{L}'_{L1'(G)}$ loss. The weighting of the)) loss function is determined by the parametric λ, which is as follows in Eq. (12).

$$G'^* = \text{Arg min}_{G'} \text{max}_{D'} \mathcal{L}'_{C'GAN(G',D')} + \lambda \mathcal{L}'_{L1'(G)} \qquad (12)$$



Fig. 2.   GAN architecture.

*D. BiGRU for Feature Selection and Classification*

RNNs include gated recurrent units (GRUs). It is additionally suggested to address issues like long-term memory and slopes in reverse propagation, which are comparable to LSTM. Using sequential information as input and all neuron linked in a series, recurrent neural networks (RNNs) conduct recursive in the developmental directions of patterns. Cells have the capacity to simultaneously acquire data from other cells and their own past events because to the presence of cyclic components in the hidden layer. As a result, storage and shared parameters are features of an RNN. RNN also performs better when training nonlinear features from serialized data [24]. Researchers offered LSTM that has the ability to acquire the correlation knowledge among lengthy immediate sequences of data, as a solution to the issue of

RNN gradient fading while able to not understand lengthy-term historical load attributes. GRU had been developed recently as a solution to the issue of LSTMs having excess parameters and a sluggish convergence rate [25]. The GRU is an LSTM variation that has fewer variables with greater convergence ability while yet retaining decent learning outcomes. On the inside, the GRU model is made up of updated gates and resetting gates. In contrast to LSTM, GRU substitutes an updated gate for inputs and forgetting gates, wherein the updated gate denotes the effect of the concealed layer of neurons' output data from a prior instant on their current state. The impact's degree is higher while the latest gate value is higher. The disregard level of the buried layer neuron output at the prior instant is represented by the resetting gate. A smaller amount of data is disregarded as the reset gate value increases. The hidden layer unit A can be calculated by the following Eq. (13) to Eq. (15)

$$b_t = \sigma(A_Z.[c_{t-1}, y_t]) \tag{13}$$

$$m_t = \sigma(A_m.[c_{t-1}, y_t]) \tag{14}$$

$$\tilde{c}_t = \tanh(A.[m_t * c_{t-1}, y_t]) \tag{15}$$

$$c_t = (1-b_t)* c_{t-1} + b_t * c_t \tag{11}$$

$A_Z$, $U_Z$ and U are all training parameters matrices, while $b_t$ and $m_t$ are the updating gate and resetting gate, correspondingly. Tanh is the hyperbola tangent value. The resetting gate $m_t$, the layer that is hidden, the neuron's output $c_{t-1}$, the currently inputted $y_t$ the trained parameter matrices, and U all work together to establish the candidate activation state $\tilde{c}_t$ at the present instant.

The ability of the BiGRU network to understand the connection among factors that have influenced previous and future demands and the current load makes it easier to extract the deep characteristics of load data. Fig. 3 shows BIGRU Architecture in [26] BIGRU architecture is displayed in Fig. 3



Fig. 3. BIGRU architecture.

### E. African Buffalo Optimization Algorithm

The African buffalo's skill is enhanced by its place of searching in the African Buffalo Optimization technique. It is widely used to identify buffaloes by where they are and the sounds they make ('Waa' and 'maa'). Additionally, studying aspects will help in the motion of the buffalo. The 'Waa' and 'maa' sounds are represented by the letters $wn$ and $mn$ respectively. Using the formula, cooperative productivity is clearly specified in Eq. (16).

$$mn + 1 = mn + le1(cemax - wn) + le2(cdmaxn - wn) \tag{16}$$

where, $mn$ and $wn$ stand for the nth buffalo's discovery and extraction moves, respectively (n=1, 2, 3, N). The variables for learning $le1$ are $le2$. In (1), $cemax$ is the herd's optimal fitness level and $cdmaxn$ is each buffalo's optimal fitness level is Update the place of the buffalo in ($cemax$ and $cdmaxn$) appears to be a part of the description provided for the African Buffalo Optimization algorithm. However, this statement lacks specific information or equations to clarify how the update process is performed. To provide a more accurate response, I would need additional details or equations specifying how the update process is carried out in Eq. (17)

$$wn + 1 = 2(wn + wn) \tag{17}$$

Three main components: (mn+1) the remembrance part, whereby the animals pay attention to being moved from their original location (mn). Broad memory ability is shown in their nomadic lifestyle, which is an essential tool for buffalos. The cooperative traits of buffalo are represented in the next section, $le1 (cemax - wn)$. Buffalos can trail the locations and are effective transmitters in every iteration. The last equation, $le2 (cemax - wn)$, reveals the superior intellect of buffaloes. They are able to compare their present position to their old, most productive job.

---

Algorithm 1: **Pseudocode of African Buffalo Optimization (ABO)**

---

Initialize the population of buffalos randomly;

Evaluate the fitness of each buffalo in the population;

Repeat until convergence:

for each buffalo in the population:

Select a random buffalo from the population;

if (fitness of the selected buffalo > fitness of the current buffalo):

Move towards the selected buffalo based on position update equation;

Perform boundary checks to ensure that the buffalo stays within the search space;

Evaluate the fitness of the updated buffalo;

Update the best solution found so far;

End Repeat

---

Return the best solution found;

---

## V. RESULT AND DISCUSSION

The innovative approach presented in this study for diabetic retinopathy identification leverages a hybrid Generative Adversarial Network (GAN) and Bi-directional Gated Recurrent Unit (BiGRU) model, refined through the application of the African Buffalo Optimization algorithm. By capitalizing on the GAN's proficiency in extracting intricate features from retinal images, the model achieves an enhanced capacity to discern subtle patterns indicative of diabetic retinopathy. The crucial aspect of feature extraction is addressed by the GAN, revealing concealed information essential for precise diagnosis. Subsequently, the BiGRU component effectively manages temporal relationships, facilitating comprehensive information absorption from the extracted features. The amalgamation of GAN's feature extraction and BiGRU's sequential information processing engenders a synergistic synergy, endowing the model with a profound understanding of retinal images. Furthermore, the utilization of the African Buffalo Optimization technique fine-tunes the model's parameters, optimizing its performance and resulting in an impressive accuracy rate of 98.5% in diabetic retinopathy detection. This Python-based study not only attests to the model's exceptional accuracy but also underscores its remarkable efficacy in advancing the field of diabetic retinopathy diagnosis.

### A. Performance Evaluation

Precision is the most often used approach for measuring categorization efficacy among the key assessment measures. By counting the percentage of test datasets that a classifier properly classifies, precision evaluates a classifier's accuracy. However, focusing entirely on accuracy might be constrained because it occasionally doesn't result in the best choices. Researchers have included others such as accuracy, recall, precision, and F1-score, to fully assess the classifier's

performance. These measurements are each defined as follows: Accuracy:

The accuracy of the model is assessed using confusion metrics, a widely used statistic that assesses the model's performance in classification tasks. The percentage of instances that were properly detected out of all the cases taken into consideration is measured by the accuracy metric (ACC). This parameter, which is frequently presented as a percentage, shows how precisely the classification algorithm can pinpoint the pertinent circumstances. Accuracy in this scenario is determined as Eq. (18)

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \qquad (18)$$

Wherein FP stands for False positives, TP represents for true positives, TN represents for true negatives, and FN is for false negatives.

Precision (P), which is the proportion of true positives to all positively identified instances, serves as a key assessment indicator in the present investigation. This is the proportion of persons correctly categorized as having the illness amongst every person, as Eq. (19)

$$Precision = \frac{TP}{TP+FP} \qquad (19)$$

Recall (R), which in this case stands for the proportion of true positives correctly identified by the model, is also quite important. Recall is an important consideration when evaluating the efficacy of the framework in the framework of this study in Eq. (20)

$$Recall = \frac{TP}{TP+FN} \qquad (20)$$

These criteria are essential for assessing how well a DR categorization system is working. They are used to construct the F1-score, which is defined theoretically as follows and stands for the harmonic mean of accuracy and recall:

$$F1 - score = 2 * \frac{Precision*Recall}{Precision+Recall} \qquad (21)$$

This F1-score serves as a crucial barometer of how well the system can identify people who are impacted by Drive Stages of DR Classification Results



Fig. 4. Training and Testing Accuracy

The graphical depiction for the training and validation accuracy of the proposed method, ABO based GAN-BIGRU Fig. 4 and Fig. 5 illustrates the performance of the model during the training process. The validation accuracy curve demonstrates how effectively the algorithm extends to new data, whereas the training accuracy curve often demonstrates the extent to which the model learns from training data. These curve' divergence and convergence patterns reveal information about the model's capacity to identify characteristics and generate precise predictions. In order to evaluate the model's performance and make sure that it does not over-fitting the information used for training, it is crucial to keep an eye on these efficiency curves.

A crucial visual depiction of the process of learning for the ABO based GAN-BIGRU is the process of training and validation loss. The model's capacity to minimize error during training and to adapt to new data during verification is shown by these loss curves. The validity loss curve demonstrates how well the algorithm extends whereas the training loss curve normally lowers as the system improves from the training data. By observing these curves, you can assess the simulation's capability to fitting the data without overfitting, identify convergence or divergence patterns, and determine whether adjustments are needed in the training process, such as modifying hyperparameters or adjusting the model architecture to improve its overall performance.

Fig. 6 displays Enhancing the fitness of the African Buffalo Optimization (ABO) algorithm involves a multifaceted approach, including parameter tuning, hybridization with other optimization methods, adaptive strategies, local search techniques, problem-specific customization, parallelization, thoughtful fitness function design, adjusted termination criteria, robustness enhancements, and rigorous experimental validation. By systematically applying these strategies, the ABO algorithm can be tailored to address a wide spectrum of optimization challenges and improve its capacity to converge to optimal or near-optimal solutions effectively.



Fig. 6. Fitness improvement over iterations (ABO).



Fig. 7. ROC curve for proposed ABO-based GAN-BIGRU.

Fig. 7 symbolizes the GAN-BiGRU model for diabetic retinopathy identification using the ABO (African Buffalo Optimization) Based ROC curve exhibits the algorithm's capability to discriminate between individuals with and without diabetic retinopathy across different categorization criteria. The area under the curve (AUC) value will increase as the ROC curve gets closer to the top-left corner of the plot, suggesting an algorithm with good diagnostic effectiveness and high true positive rates and low false positive rates. The properties of the dataset and the model projections will determine the precise curve shape and AUC value.



Fig. 5. Training and testing loss.

Fig. 8.   An illustration of interestingness during instruction.

Fig. 8 illustrates interestingness during instruction According to the statement; interest is maintained throughout the training period. The interestingness value stays non-zero while the model learns from the training data and modifies its

parameters, suggesting continuous updates that support the model's learning process. This continuous attention is probably related to how the model has adjusted to the subtleties and patterns found in the training set.

Table I compares the suggested strategy with several current methods in terms of performance measures. This tabled data provides a summary of the major assessment criteria where the suggested technique outperforms its predecessors. It is a useful resource for comprehending the better performance and efficacy of the suggested strategy in the particular application or study topic.

TABLE I.        ASSESSMENT OF PERFORMANCE METRICS OF SUGGESTED METHOD WITH FURTHER EXISTING APPROACHES

| Methods | Accuracy (%) | Precision (%) | Recall (%) | F1-score (%) |
|---|---|---|---|---|
| Alex Net[27] | 97.9 | 96.23 | 95.42 | 795.82 |
| Random Forest [28] | 94.1 | 97.6 | 94.3 | 95.9 |
| VGG-NIN[15] | 94.20 | 90.0 | 98.0 | 94.0 |
| Proposed ABO Based GAN-BIGRU | 98.5 | 98.54 | 96.3 | 97.8 |



Fig. 9.   Visual representation comparing the proposed method with existing approaches.

Fig. 9 shows comparison of proposed method with existing methods. Alex Net (achieved an accuracy of 97.9%, with a precision of 96.23%, a recall of 95.42%, and an F1-score of 95.82%. Random Forest exhibited an accuracy of 94.1%, with a high precision of 97.6%, a recall of 94.3%, and an F1-score of 95.9%.VGG-NIN showed an accuracy of 94.20%, with a precision of 90.0%, a notably high recall of 98.0%, and an F1-

score of 94.0%. The proposed ABO Based GAN-BIGRU model achieved an impressive accuracy of 98.5%. It also demonstrated high precision at 98.54%, a recall of 96.3%, and a remarkable F1-score of 97.8%. These results highlight the varying performance of different classification models for diabetic retinopathy, with the "Proposed ABO Based GAN-BIGRU" model showing the highest overall performance across all metrics.

TABLE II. RETINAL IMAGE DATASETS AND THEIR ATTRIBUTES

| Dataset Name | DR Lesion Annotation | Vessel Segmentation | Resolution | Number of Images |
|---|---|---|---|---|
| DRIVE | (NO DR) | YES | 896×896 | 10 train + 10 test |
| Kaggle | Only severity levels | YES | 1281×1281 | 36.1k train + 54.6k test |
| IDRiD | Pixel-wise lesion segmentation & severity levels | NO | 4289×2849 | IDRiD Pixel-wise lesion segmentation No 4289×2849 |
| Retinal-Lesions | Pixel-wise lesion segmentation | NO | 583×575 | 338 train + 1257 test |
| FGADR | Lesion annotation in circle & severity level | YES | 1280×1280 | 500 train + 1343 test |

Retinal image datasets and attributes are shown in Table II Several retinal image datasets with diverse attributes have been instrumental in advancing diabetic retinopathy research. The DRIVE dataset, characterized by 10 training and 10 testing images, focuses on vessel segmentation with a resolution of 896x896 pixels. The Kaggle dataset, comprising 36.1k training and 54.6k testing images, annotates only severity levels and includes vessel segmentation at a resolution of 1281x1281 pixels. IDRiD offers pixel-wise lesion segmentation and severity levels with a substantial resolution of 4289x2849 pixels. The Retinal-Lesions dataset features pixel-wise lesion segmentation in 338 training and 1257 testing images with a resolution of 583x575 pixels. Finally, the FGADR dataset provides lesion annotation in circles, along with severity levels, and incorporates vessel segmentation in 500 training and 1343 testing images at a resolution of 1280x1280 pixels. These datasets play a crucial role in fostering the development and evaluation of diabetic retinopathy detection models, each offering unique challenges and opportunities for research and advancement in the field.

## B. Discussions

The hybrid Generative Adversarial Network (GAN) and Bi-directional Gated Recurrent Unit (BiGRU) model, fine-tuned with the African Buffalo Optimization algorithm, represents a notable advancement in diabetic retinopathy detection. By combining GAN's proficiency in extracting intricate features from retinal images with BiGRU's effective handling of temporal relationships, the model achieves an impressive 98.5% accuracy rate. The incorporation of the African Buffalo Optimization algorithm further optimizes the model's parameters, showcasing the potential of bio-inspired optimization techniques in enhancing the performance of deep learning models for medical image analysis. This integrated approach not only underscores the significance of feature extraction and temporal considerations in diabetic retinopathy diagnosis but also highlights the promising synergy achievable through the convergence of diverse neural network architectures and optimization strategies.

## VI. CONCLUSION AND FUTURE WORK

In conclusion, the hybrid GAN-BiGRU model, fine-tuned with the African Buffalo Optimization algorithm, presents a robust solution for diabetic retinopathy detection, achieving an impressive accuracy rate of 98.5%. This study underscores the effectiveness of combining advanced feature extraction capabilities with temporal information processing, showcasing the synergy between GAN and BiGRU architectures. The integration of the African Buffalo Optimization algorithm further refines the model's parameters, emphasizing the potential of bio-inspired optimization in enhancing diagnostic accuracy. For future work, exploring the model's generalizability across diverse populations and investigating its applicability to real-world clinical settings would be valuable. Additionally, continued efforts to interpret the model's decisions and address potential biases in the dataset could contribute to its clinical adoption. Further refinement and validation through large-scale multi-center studies could solidify the model's potential as a valuable tool in early diabetic retinopathy diagnosis.

## REFERENCES

[1] Y. Xu et al., "Global burden and gender disparity of vision loss associated with diabetes retinopathy," Acta Ophthalmol. (Copenh.), vol. 99, no. 4, pp. 431–440, 2021.

[2] K. Suriyasekeran, S. Santhanamahalingam, and M. Duraisamy, "Algorithms for diagnosis of diabetic retinopathy and diabetic macula edema-a review," Diabetes Res. Clin. Pract. Vol. 4, pp. 357–373, 2021.

[3] S. Joseph, R. P. Rajan, B. Sundar, S. Venkatachalam, J. H. Kempen, and R. Kim, "Validation of diagnostic accuracy of retinal image grading by trained non-ophthalmologist grader for detecting diabetic retinopathy and diabetic macular edema," Eye, vol. 37, no. 8, pp. 1577–1582, 2023.

[4] Y. Sun, "The neural network of one-dimensional convolution-an example of the diagnosis of diabetic retinopathy," IEEE Access, vol. 7, pp. 69657–69666, 2019.

[5] R. Manne and S. C. Kantheti, "Application of artificial intelligence in healthcare: chances and challenges," Curr. J. Appl. Sci. Technol., vol. 40, no. 6, pp. 78–89, 2021.

[6] J. V. Forrester, L. Kuffova, and M. Delibegovic, "The role of inflammation in diabetic retinopathy," Front. Immunol., vol. 11, p. 583687, 2020.

[7] P. J. Snyder et al., "Retinal imaging in Alzheimer's and neurodegenerative diseases," Alzheimers Dement., vol. 17, no. 1, pp. 103–111, 2021.

[8] C. Singh, M. Mallesha, M. Vijayaragavan, J. SURESHBABU, and others, "IoT based secured healthcare using 6g technology and deep learning techniques," J. Pharm. Negat. Results, pp. 462–472, 2022.

[9] D. A. Antonetti, P. S. Silva, and A. W. Stitt, "Current understanding of the molecular and cellular pathology of diabetic retinopathy," Nat. Rev. Endocrinol., vol. 17, no. 4, pp. 195–206, 2021.

[10] P. Uppamma, S. Bhattacharya, and others, "Deep Learning and Medical Image Processing Techniques for Diabetic Retinopathy: A Survey of Applications, Challenges, and Future Trends," J. Healthc. Eng., vol. 2023, 2023.

[11] A. Grzybowski and P. Brona, "Analysis and comparison of two artificial intelligence diabetic retinopathy screening algorithms in a pilot study: IDx-DR and RetinaLyze," J. Clin. Med., vol. 10, no. 11, p. 2352, 2021.

[12] M. Sharif, J. H. Shah, and others, "Automatic screening of retinal lesions for grading diabetic retinopathy," Int Arab J Inf Technol, vol. 16, no. 4, pp. 766–774, 2019.

[13] A. Tariq, A. Y. Gill, and H. K. Hussain, "Evaluating the Potential of Artificial Intelligence in Orthopedic Surgery for Value-based Healthcare," Int. J. Multidiscip. Sci. Arts, vol. 2, no. 1, pp. 27–35, 2023.

[15] P. Ebin and P. Ranjana, "An approach using transfer learning to disclose diabetic retinopathy in early stage," in 2020 International Conference on Futuristic Technologies in Control Systems & Renewable Energy (ICFCR), IEEE, 2020, pp. 1–4.

[16] Z. Khan et al., "Diabetic retinopathy detection using VGG-NIN a deep learning architecture," IEEE Access, vol. 9, pp. 61408–61416, 2021.

[17] L. Qiao, Y. Zhu, and H. Zhou, "Diabetic retinopathy detection using prognosis of microaneurysm and early diagnosis system for non-proliferative diabetic retinopathy based on deep learning algorithms," IEEE Access, vol. 8, pp. 104292–104302, 2020.

[18] M. M. Farag, M. Fouad, and A. T. Abdel-Hamid, "Automatic severity classification of diabetic retinopathy based on densenet and convolutional block attention module," IEEE Access, vol. 10, pp. 38299–38308, 2022.

[19] D. J. Hemanth, O. Deperlioglu, and U. Kose, "An enhanced diabetic retinopathy detection and classification approach using deep convolutional neural network," Neural Comput. Appl., vol. 32, pp. 707–721, 2020.

[20] R. R. Maaliw et al., "An Enhanced Segmentation and Deep Learning Architecture for Early Diabetic Retinopathy Detection," in 2023 IEEE 13th Annual Computing and Communication Workshop and Conference (CCWC), IEEE, 2023, pp. 0168–0175.

[21] W. Chen, B. Yang, J. Li, and J. Wang, "An approach to detecting diabetic retinopathy based on integrated shallow convolutional neural networks," IEEE Access, vol. 8, pp. 178552–178562, 2020.

[22] S. Qummar et al., "A deep learning ensemble approach for diabetic retinopathy detection," Ieee Access, vol. 7, pp. 150530–150539, 2019.

[23] Y. Chen et al., "Automated cone photoreceptor cell segmentation and identification in adaptive optics scanning laser ophthalmoscope images using morphological processing and watershed algorithm," IEEE Access, vol. 8, pp. 105786–105792, 2020.

[24] A. Sebastian, O. Elharrouss, S. Al-Maadeed, and N. Almaadeed, "GAN-Based Approach for Diabetic Retinopathy Retinal Vasculature Segmentation," Bioengineering, vol. 11, no. 1, Art. no. 1, Jan. 2024, doi: 10.3390/bioengineering11010004.

[25] M. B. Darici and G. Yigit, "Improving Diabetic Retinopathy Detection Using Patchwise CNN with biGRU Model," in 2023 8th International Conference on Computer Science and Engineering (UBMK), IEEE, 2023, pp. 1–5.

[26] Q. Yu, Z. Wang, and K. Jiang, "Research on text classification based on bert-bigru model," in Journal of Physics: Conference Series, IOP Publishing, 2021, p. 012019.

[27] H. Qiu, C. Fan, J. Yao, and X. Ye, "Chinese Microblog Sentiment Detection Based on CNN-BiGRU and Multihead Attention Mechanism," Sci. Program., vol. 2020, pp. 1–13, Oct. 2020, doi: 10.1155/2020/8865983.

[28] N. Khalifa, M. Loey, M. Taha, and H. Mohamed, "Deep Transfer Learning Models for Medical Diabetic Retinopathy Detection," Acta Inform. Medica, vol. 27, no. 5, p. 327, 2019, doi: 10.5455/aim.2019.27.327-332.

[29] N. P. Tigga and S. Garg, "Prediction of type 2 diabetes using machine learning classification methods," Procedia Comput. Sci., vol. 167, pp. 706–716, 2020.

# The Application of Artificial Intelligence Technology in Ideological and Political Education

Chao Xu[1]*, Lin Wu[2]

School of Marxism, Xinyang University, Xinyang, Henan 464000, China[1]
School of Education, Xinyang University, Xinyang, Henan 464000, China[2]

*Abstract*—As for many schools, artificial intelligence will be more than a practical background; it is also a technical tool and an opportunity for development. Artificial intelligence's in-depth integration and standardization can inject new technological momentum into effectively identifying educational objects' ideological dynamics, improving educational content's accuracy, and expanding the spatial dimension. It has become one and such an inevitable trend of innovation and development. However, there are also many potential risks and practical problems at the value premise, technical limits, and specific operation level, such as privacy protection and ideological security risks, the loss of educational subjectivity, the digitization of educational relations, and the lack of specialized talents. Therefore, it is necessary to look at the technical momentum and potential risks of artificial intelligence dialectically, promote the rationality of educational value, strengthen technical supervision, forge an intelligent education team, reasonably define the integration boundary and application scope of artificial intelligence, and combine the main initiative of human beings with the intelligence of machines. It combines strengths, actively explores the path of coexistence and co-prosperity between education and technology, and consciously constructs an intelligent form of them.

*Keywords—Artificial intelligence; ideological and political education; wisdom development; semantic understanding and emotional analysis*

## I. INTRODUCTION

As the frontier one, to usher the ones in the intelligent era"[1], the wide fields of society will be prompting another profound change in education, which has profoundly impacted various educational practice activities [2]. In the face of the superimposed development of informatization, the Internet, blockchain, and artificial intelligence (AI) technology, He has repeatedly stressed that in education and other fields, innovate the intelligent service system"[3].

The dynamic changes in the attention and interest of educational objects track the interactive data between educators and educational objects on time and investigate the educational objects through data change trajectories. The state of thinking and behavior accurately grasp the individual needs of educational objects and provides reliable support for formulating scientific and effective educational programs. For artificial intelligence, the processing technology transcends traditional statistical analysis methods, can collect and analyze all data samples, and can integrate qualitative and quantitative data, historical and current data.

It is necessary to think deeply about how to update one iteratively, follow the trend, continuously strengthen the deep integration with artificial intelligence, and use the power of artificial intelligence to promote intelligent transformation and upgrade them to provide a high-quality future society [4]. The development of socialist modernization and the comprehensive construction of socialist modernization cultivate intelligent compound talents [5].

As a powerful scientific and technological force that subverts traditional educational concepts and shapes the future education form, its accelerated iteration and all-around integration have opened a new chapter in innovative development and have classified such objects and contents [6]. Quality improvement in distribution, space expansion, and discourse expression provides strong technical support [7].

"If Communist Party members want to do propaganda, they must look at the objects." The prerequisites are effective classification of objects, and accurate identification of object needs [8]. Subject to factors such as ideological complexity and technical limitations, traditional ones usually distinguish objects based on criteria such as field, age, major, and class, which makes it difficult to reflect the particularity, complexity, and dynamics of education [9]. For artificial intelligence, "the application of technology allows us to record them more completely, and provides convenience for using other research methods to understand people's thoughts."

People's thoughts and concepts, emotional orientations, hobbies, and daily behaviors in the real world may be mapped to some digital existence, and people's thoughts and behaviors thus become measurable, recordable, predictable, and cluster-analyzable." Intelligent data becomes the second body of human beings [10]." "Once the relevant data was fully grasped, it is possible to completely predict and grasp individuals' behavioral tendencies and dispositions using digital drawing [11]." The data portrait of people's thoughts and behaviors opens new possibilities for people to understand themselves and profoundly changes them [12]. Their ecology has been realized, and the revolutionary leap for classification technology of ideological and political education objects has been established, thus establishing the internal relationship between the two [13].

For artificial intelligence, the processing technology transcends such traditional statistical analysis methods, can collect and analyze all data samples, and can integrate qualitative data and quantitative data, historical data and current data, individual data and overall data [14]., to more

completely outline the appearance of human thought and behavior [15].

First, objectively present the ideological and behavioral state of educational objects. Under the condition of AI empowerment, colleges and universities can, within the scope of laws and ethics, comprehensively collect data traces such as social preferences, browsing preferences, and value orientations of educational objects in their daily lives and clarify their different roles and responsibilities in the online world [16]. Artificial intelligence (AI) has penetrated into various fields, bringing unprecedented convenience to our lives and work. In the field of education, especially in ideological and political education in universities, the introduction of AI is changing the traditional education model, providing a more personalized and efficient learning experience [17]. AI can provide real-time learning guidance and troubleshooting for students through learning and analyzing a large amount of data, improving their learning effectiveness. AI can simulate human emotional communication, establish emotional connections with students, and enhance their learning motivation and satisfaction [18].

The second is to dynamically capture the changes in the thinking and behavior of educational objects [19]. With the help of technologies such as attention recognition, dynamic capture, and data association analysis, colleges and universities can correctly observe the dynamic changes in the attention and interest of educational objects, track the interactive data between educators and educational objects on time, and investigate the educational objects through data change trajectories [20]. The state of thinking and behavior accurately grasp the individual needs of educational objects and provides reliable support for formulating scientific and effective educational programs. For artificial intelligence, the processing technology transcends traditional statistical analysis methods, can collect and analyze all data samples, and can integrate qualitative and quantitative data, historical and current data [21].

Third, educational object classification standards are more scientific, and the presentation method is more intuitive. Under the background of artificial intelligence technology, it has become a reality to classify educational objects based on the standard of "difference in ideas"[22]. This standard helps classify educational objects by highlighting the individual's ideas, behavioral orientation, emotional attitude, hobbies, and differences. More importantly, artificial intelligence technology can turn college students' dynamically changing ideological behaviors, emotional concepts, and other vague elements into clear, accurate data and present different circles and types. This provides a technical prerequisite for them to carry out various activities for batches, circles, and methods [23].

It can simulate any other prediction function p(x) of the algorithm. This is not possible with any other machine learning algorithm [24]. One is the imitation theory. The theory of imitation mainly believes that artificial intelligence explores the thinking process of the human brain, mainly simulates

human intelligence and studies the science of extending human mental work to some physical device. Artificial intelligence is the ability of machines (computers) to perform some complex functions related to human intelligence (such as judgment, pattern recognition, understanding, learning, decision-making, planning, problem-solving, etc.). Artificial intelligence's task is to replace the human brain and manual labor with machines partially. For example, some scholars think that the purpose of AI is to make the computer, the machine, think like a human. Artificial intelligence is a technology that simulates the realization of human thinking. Its main purpose is to give robots the unique ability to see, hear, speak, and abstract thinking in the brain. It is especially reflected in thinking activities such as judgment, reasoning, proof, identification, learning, and problem-solving. In general, it is a combination of knowledge and thinking. This statement is widely used in academic circles and imitation of human intelligence, imitating human thinking and action patterns [25].

The main contributions of this study are as follows:

*1)* Artificial intelligence can provide personalized educational programs for students by analyzing their learning habits, interests, and career plans. This educational approach helps to stimulate students' interest in learning, improve their learning outcomes, and also helps to enhance the pertinence and effectiveness of ideological and political education.

*2)* Artificial intelligence technology can assist teachers in teaching management, course design, and teaching evaluation. This helps to reduce the workload of teachers, improve teaching efficiency, and also enhance the quality of ideological and political education.

*3)* By building an intelligent learning platform, students can learn anytime and anywhere. This learning method helps to break the limitations of traditional classrooms and provide students with more flexible and convenient learning methods.

Section I first analyzes the advantages of AI empowerment. Universities can dynamically capture changes in the thinking and behavior of educational objects, and comprehensively collect their preferences. Section II introduces the relevant concepts and theoretical foundations. Section III utilized RCGA, which effectively addressed multiple domains. The focus is on the ability of wavelet functions to process synthetic data, as well as the classification of functions, hyperbolic tangent functions, and threshold linear functions. Section IV conducted a survey and analysis on the value demands of empowering college students with artificial intelligence in ideological and political education. Decompose the learning task of fuzzy cognitive maps based on their sparsity. This algorithm has overcome the shortcomings of existing methods and achieved breakthroughs of multiple orders of magnitude. Section V summarizes the entire text, and it is necessary to dialectically examine the technological momentum and potential risks of artificial intelligence, promote the rationality of educational value, strengthen technical supervision, build an intelligent education team, and reasonably define the integration boundary and application scope of artificial intelligence.

## II. RELATED CONCEPTS AND THEORETICAL BASIS

Its related inventions and applications have begun to cover our lives, and intelligent life has become within reach [26]. The research basis of the present work is the artificial intelligence environment and its use for the problems caused by traditional political and ideological education of college students, such as the inability to large-scale individualization in the process and for this chapter, a discussion will be held in conjunction with basic theoretical issues related to artificial intelligence, including the basic concepts and characteristics of artificial intelligence and the basic concepts, the new features of artificial intelligence and its application for education. Rather than optimizing some parameters in existing models, such as polynomial curves and nodal systems, the approach employed by neural networks is a specific perspective on data modeling that does not seek to exploit any independent system fully but directly approximates data functions. The neural network architectures are so familiar to model merely ideas. With the power of neural networks and continued research into the bottomless field of deep learning, data—whether video, sound, epidemiological data or anything in between—can be modeled; neural networks are indeed algorithms of algorithms [27].

### A. Artificial Intelligence

A concept is a form of thinking that reflects the essential properties of an objective object. Understanding the concept of AI, especially the concept of intelligence, can effectively reflect the essential characteristics of artificial intelligence and is helpful. In academia, intelligent science and technology are the core of research in nature, humanities, and social sciences effects and changes [28].

This is compared with the concept of intelligence in the context of applied science, which was proposed as early as the 17AC and was first proposed in 1956. Still, in the short decades after it was proposed, artificial intelligence developed rapidly. Robot vision, intelligent robots, robot planning, and other categories include intelligent control in modern control theory, fuzzy mathematics, and fuzzy control theory in modern control theory [29].

There are many opinions in the academic circle about the definition of artificial intelligence. Among them are mainly divided into:

One is the imitation theory. The theory of imitation mainly believes that artificial intelligence explores the thinking process of the human brain, mainly simulates human intelligence and studies the science of extending human mental work to some physical device. Artificial intelligence is the ability of machines (computers) to perform some complex functions related to human intelligence (such as judgment, pattern recognition, understanding, learning, decision-making, planning, problem-solving, etc.). Artificial intelligence's task is to replace the human brain and manual labor with machines partially. For example, some scholars think that the purpose of AI is to make the computer, the machine, think like a human. Artificial intelligence is a technology that simulates the realization of human thinking. Its main purpose is to give robots the unique ability to see, hear, speak, and abstract thinking in the brain. It is especially reflected in thinking activities such as judgment, reasoning, proof, identification, learning, and problem-solving. In general, it is a combination of knowledge and thinking. This statement is widely used in academic circles and imitates human intelligence, thinking, and action patterns [30].

The second is the expansion theory. The expansion theory of the concept of artificial intelligence mainly believes that artificial intelligence is not only an imitation of human intelligence but should play a role in expanding human intelligence based on imitating human thinking and behavior and finally achieving the goal of enhancing human intelligence purpose.

Third, comprehensively. The comprehensive theory of artificial intelligence concept mainly refers to artificial intelligence, including all the terms of many sub-fields, involving an extensive range of applications. In the view of scholars who support the comprehensive theory, artificial intelligence does not only refer to robots that can imitate human thinking and behavior. They believe that artificial intelligence includes technologies such as image recognition, video recognition, semantic understanding, and sentiment analysis. It is a general term for specific technologies but does not regard artificial intelligence as a general ability.

This paper is considered the manifestation of the prosperity of deep learning algorithms based on big data, and it is not equivalent to the "general artificial intelligence" that tried to restore human intelligence and behavior in the form of robots in the past. It is not that only robots are artificial intelligence, as most people think. Artificial intelligence should be a comprehensive technology with six major technical directions: big data, statistical analysis, natural language processing, speech processing, planning and decision-making systems, and computer vision. One or several technical directions can be called artificial intelligence. Artificial intelligence can perceive, analyze, understand, think, decide, and interact. The ones (cloud computing) are the three cornerstones of today's artificial intelligence technology, of which the main technologies of artificial intelligence are deep learning and big data analysis. This paper is to conduct further research on this basis.

### B. The Ideological and Political Education of College Students

Since the birth of New China in 1949, generations of outstanding college students have played a pivotal and historic role in gradually building our country into a prosperous, democratic, civilized, and harmonious modern socialist country. As pointed out in the one. On September 30, 2013, the ninth collective study of the Political Bureau of the Eighteenth Central Committee found that it is an important source of talent for socialist construction. As the primary position in the education of college students, it has great strategic significance in cultivating them into qualified construction. Logical reasoning is one of the most enduring for AI research. It is important to find ways to focus only on relevant facts in a large database, be on the lookout for credible proofs, and revise them as new information emerges. Finding a proof or disprove of a speculative theorem in mathematics is indeed an intelligent task.

*1) Group characteristics of college students:* They are mainly born after "95" and "00". Compared with the "post-90s", the group has more obvious group characteristics. The ideological and political education of college students must be targeted according to the group features of current college students, the pursuit of individualized value. "Post-00" college students grew up in an era that advocated independence and advocated that contemporary college students should have independence and autonomy. Therefore, when faced with choices, the new generation of college students who grew up in a material-rich environment tend to ignore material pursuits and pay more attention to personal emotional experience and the realization of self-worth. Without large-scale personalized teaching, it is impossible to meet the personalized pursuit of every college student. The characteristics of college students' value pursuit in the new era have brought challenges to the development of college students' ideological and political education.

Second, the network behavior is diversified. "Post-00" college students were born and grew up in the Internet age, and their daily lives and studies reflect "Internet thinking." As the "indigenous people" of the Internet, their daily communication is full of Internet terms; most of their communication methods use WeChat, QQ, and other software; the proportion of their online consumption far exceeds that of brick-and-mortar stores. The Internet has become the "residence place" of college students after "00". At this time, ideological and political education in college students should conform to the background of the Internet era, use big data, artificial intelligence technology, etc., and its relevance and effectiveness. Robots that can imitate human thinking and behavior. It includes technologies such as image recognition, video recognition, semantic understanding, and sentiment analysis. It is a general term for specific technologies but does not regard artificial intelligence as a general ability.

Therefore, the students should conform to the development of the times, and under the clear background of the times, macroscopically grasp the character characteristics of the entire college student group and carry out one and political education based on this. The group of college students who are mainly "post-00s" happens to be in the social background of the rise of artificial intelligence. It is vital to think about how to change the ideology and politics and whether to change or stick to it.

*2) Main features:* It has entered a new era, and it is required to firmly grasp the leadership of ideological work, cultivate and practice socialist core values, and strengthen ideological and moral construction. College students should also have a new form. As the main direction of technological innovation, the integration of AI with ideological and political education in college students is imperative, and it is conducive to the "intelligence enhancement" for artificial intelligence; learners should have new ideas, new forms, new models and new paths, in Fig. 1:



Fig. 1.   Smart teaching mode of ideological and political education for college students.

The theory of human beings is the core theory of Marxist theory. Marx believes it mainly refers to such labor ability, social relations, the all-round development and great satisfaction of human needs, and the freedom of human personality. The comprehensive for college students lies in the comprehensive development of comprehensive quality, specifically the comprehensive, coordinated, and sustainable one.

Under the new historical conditions, there are differences in the individual ones. To achieve the all-round development of every college student, relying only on the existing forms of ideological and political education for college students is far from achieving the goal of promoting the all-around development in the general environment, mainly carried out in the form of ideological and political theory courses. The forms and contents of ideological and political education are gradually enriched with the reform. It is still impossible to consider every college student, and it is impossible to target education according to the personality characteristics of college students. With the blessing of artificial intelligence, they can conduct big data analysis based on the data intelligently sensed and excavated by artificial intelligence. By building a visual personality model of college students, they know the ideological development characteristics of college students and realize precise and personalized education., to achieve the all-round development of college students. College students have also changed from passively to being able to do independently and completing the basic required learning content to achieve free development.

Therefore, in the future, artificial intelligence will enable the emergence of a new concept of precise and personalized education, which can realize high-precision personalized education through big data and make large-scale personalized education a reality. Not only small-scale personalized education can be achieved.

## III. RELATED TECHNOLOGIES

The RCGA-FCM algorithm automatically generates FCM by using historical data without human intervention. This method can provide a fully automated solution for learning FCMs. In addition, the algorithm also compared the size and connection density of FCM to select the most effective design environment. This chapter uses RCGA, and it has effectively dealt with multiple domains. Then, the chapter needs it, and this chapter focuses on the ability of wavelet functions to deal with synthetic data and classification problems for functions, hyperbolic tangent functions, and threshold linear functions rather than one. When learning such FCM, it will be

$$W = \left[ w_{11}, w_{12}, \dots, w_{1N_n}, w_{21}, \dots, w_{2N_n}, \dots, w_{N_nN_n} \right] \quad (1)$$

The wavelet function-based FCM one is designed. The one which is an approximation of such one:

$$F = \cfrac{1}{\left( \cfrac{\beta}{N(T-1)} \sum_{t=1}^{T} \sum_{i=1}^{N} (C_i^*(t) - C_i(t))^2 + 1 \right)} \quad (2)$$

### A. Wavelet Fuzzy Cognitive Map

The wavelet function-based FCM one is designed. The one which is an approximation of such one. A wavelet function consists of the following:

$$L^2(R) = \left\{ x(t) :_R \int |x(t)|^2 dt < \infty \right\} \quad (3)$$

The mother wavelet $\psi(x) \in L2(R)$ must satisfy the condition

$$L^2(R) = \left\{ x(t) :_R \int |x(t)|^2 dt < \infty \right\} \quad (4)$$

where, it represents the wavelet family achieved by the parent wavelet $\psi(x)$ by dilation and translation. It will be like this:

$$\psi_{a,b}(x) = |b|^{-\frac{1}{2}} \psi \left( \frac{x-a}{b} \right) \quad (5)$$

where, $\psi a,b(x)$ represents the wavelet family achieved by the parent wavelet $\psi(x)$ by dilation and translation. Reference [17] carried out a detailed analysis of wavelets. Based on Eq. (3) and Eq. (5), the kinetic equation of WFCM is:

$$C_i(t+1) = \psi_{a,b} \left( \sum_{j=1}^{N_n} w_{ji} C_j(t) \right) \quad (6)$$

Such artificial data and it then generates response sequences from each initial state vector by implementing step c) until T iterations are reached.

$$\psi_{a,b}(x) = \left( 1 - \left( \frac{x-a}{b} \right)^2 \right) e^{\left( -\frac{(x-a)^2}{2b^2} \right)} \quad (7)$$

where, a and b are parameters in WFCM. Since Eq. (7) is easier to handle than Eq. (5), the factor of |b|-0.5 is removed in this chapter.

### B. Manual Data Regulation of Network Data

The algorithm will have the possibility of accurately learning fuzzy cognitive maps. The algorithm breaks through the shortcomings of existing methods and increases the learning scale of FCM from 40 to 1000 nodes, achieving breakthroughs of multiple orders of magnitude.

This chapter uses the following steps:

*1)* Typically, the target FCM is set to 20% or 40%, and each non-zero weight is randomly produced in [-1,1]. Note that the absolute value of each non-zero weight and other weights are set to 0;

*2)* Randomly generate initial state values belonging to [0, 1] and assign them to each node;

*3)* Use Eq. (2) and Eq. (3) to generate the next state value of each node;

*4)* Such artificial data contains S response sequences, each with T iterations, then generate response sequences from each initial state vector by implementing step c) until T iterations are reached.

The first metric is Data Error, which the available ones:

$$\text{Data Error} = \frac{1}{NS(T-1)} \sum_{i=1}^{N} \sum_{k=1}^{S} \sum_{t=1}^{T} \left( C_i^{k^*}(t) - C_i^k(t) \right)^2 \quad (8)$$

Ns (Ns=10) that differ from the initial state vectors applied in the learning process.

$$\text{Out\_of\_Sample\_Error}$$
$$= \frac{1}{NN_s(T-1)} \sum_{k=1}^{N_s} \sum_{t=1}^{T} \sum_{i=1}^{N} \left| C_i^{k^*}(t) \right. \quad (9)$$
$$\left. - C_i^k(t) \right|$$

Model_Error is to compare the weight of the learned FCM with the weight of the original FCM.

$$\text{Model\_Error} = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} \left| w_{ij} - w'_{ij} \right| \quad (10)$$

where, wij' is the weight from node i to node j in the candidate FCM.

There will be no connection between nodes; otherwise, there is a connection. Accordingly, SS_Mean is computed as:

$$\text{SS\_Mean} = \frac{2 \times \text{Specificity} \times \text{Sensitivity}}{\text{Specificity} + \text{Sensitivity}} \quad (11)$$

In

$$\text{Specificity} = \frac{TP}{TP + FN} \quad (12)$$

TP is to assess the classification ability of a classifier quantitatively; Two observers' agreement is measured using Cohen's kappa,

$$\text{kappa} = \frac{p_o - p_e}{1 - p_e} \quad (13)$$

Such observational data hypothesized such observational data to calculate the probability that it is like:

To compare the effects of one:

$$\psi(x) = e^{\left( -\frac{(x-t)^2}{2d^2} \right)} \quad (14)$$

where, d and t denote the parameters of the Gaussian function. In the present experiment, t=3, d=30. A linear function can be defined as follows:

$$\psi(x) = |x| \quad (15)$$

In the paper, Fig. 2 shows the image segmentation performance of an improved fuzzy clustering WFCM algorithm on artificial data generated by Gaussian functions. Fig. 2 shows the performance of using the RCGA-FCM evolutionary algorithm for image segmentation. This includes evaluation indicators such as accuracy, boundary clarity, and noise suppression. This image segmentation process is achieved through an evolutionary algorithm called RCGA-FCM. It uses artificial data generated by Gaussian functions, which simulate pixel intensity or grayscale values in actual images. Allow data points to belong to multiple clusters, each with a membership degree. This method is particularly suitable for image segmentation as it can better handle the boundaries and noise between pixels. RCGA will be computed and formed in the Fig. 2:



Fig. 2. Performance of WFCM on artificial data generated by Gaussian functions.

## IV. EXPERIMENTAL RESULTS AND ANALYSIS

### A. *Investigation and Analysis of the Value Appeal of AI Empowering College Students' Ideological and Political Education*

It also become a development trend. In promoting the application of AI, the subjective feelings of college students cannot be ignored. In this regard, questionnaires and other forms should be used to objectively and rationally analyze college students' attitudes, opinions, and value demands on the ideological and political education empowered by AI robots that can imitate human thinking and behavior. Technologies include image recognition, video recognition, semantic understanding, and sentiment analysis. It is a general term for specific technologies but does not regard artificial intelligence as a general ability.

According to the survey, most ones, like problems in understanding knowledge points and insufficient sense of acquisition, exist only when there are difficulties as they want to change. College students anticipate it to be more aligned with their personal experiences than a theory elevated to the point where they feel no real benefit.

When investigating the question "What do you think are the difficulties encountered for theory courses" (see Fig. 3), 732 students chose "difficulty in remembering knowledge points," accounting for 71.07%; 585 students chose the item "difficult to understand," accounting for 56.8% of the respondents; 517 students chose the item "difficult to review," accounting for 50.19% of the respondents; 354 Students. The item "unbalanced course resources" was selected, accounting for 34.37% of the studied population; 350 students chose the item "lack of acquisition," accounting for 33.98% of the population. It can be seen that the difficulty of remembering knowledge points is a common difficulty faced by college students in ideological and political education, regardless of gender, education, and professional background. According to the survey results, fairness in education is the biggest challenge faced by college students, followed by the question of whether they don't feel like they've gained anything. According to the survey, most ones, like problems in understanding knowledge points and insufficient sense of acquisition, exist only when

there are difficulties as they want to change. College students expect it to be close to their authentic personal life rather than a high above theory so that they have no real sense of gain.

When asking college students, "Do you think artificial intelligence has brought about a change in ideas?" (see Fig. 4), 676 students chose to diversify teaching methods, accounting for 65.63% of the surveyed; 584 573 students chose personalized teaching content, accounting for 55.63% of the respondents; 541 students chose intelligent resource search, accounting for 55.63% of the respondents 52.52% of the total number of students; 484 students chose the scientific evaluation method, accounting for 46.99% of the surveyed number; 454 students chose the dynamic teaching process, accounting for 44.08% of the surveyed number; 421 students chose teaching The environment is intelligent, accounting for 40.87% of the studied population; 243 students chose to reduce the burden of the classroom, accounting for 23.59% of the surveyed population.

Diversification, precision, and personalization—the goals and outcomes that political education seeks to accomplish—are the most often mentioned keywords when most students talk about the concept change, according to the survey results above. Their workload from school will not decrease with the arrival of artificial intelligence, and issues with assessments, like forgetting knowledge points, will persist.

To what extent the respondents know about artificial intelligence technology and its application, the statistical results in Fig. 5, but know nothing about its application" is the highest, indicating that most students currently cognition of artificial intelligence remains in the booking stage, but there is a lack of understanding of how artificial intelligence is applied; among the teachers, the proportion of "understanding some artificial intelligence technologies and partially understanding their applications" for the courses has a certain degree of mastery from theory to practice, but the degree is not too deep; the proportion of practitioners who "have a deep understanding of artificial intelligence technology, but have a partial understanding of its application" The highest, indicating that the practitioners group has made certain breakthroughs in the cognition and practice.



Fig. 3. Difficulties of college students in learning ideological and political theory courses.

Fig. 4. The concept change brought by artificial intelligence to the ideological and political education of college students.

Fig. 5. How much is known about AI technology and its applications.

Regarding which technologies of artificial intelligence are most easily applied to art and design education, the statistical ones in Fig. 6 show that the top three are natural language understanding (AI translation, question answering system, etc.), computer vision (image understanding, 3D vision, dynamic vision, etc.), machine vision, etc.), biometric ones and the proportions are relatively close, all around 30%, indicating that from the perspectives of students, teachers, and practitioners, the current three This AI art and design education. According to the survey, most ones, like problems in understanding knowledge points and insufficient sense of acquisition, exist only when there are difficulties as they want to change. College students expect it to be close to their authentic personal life rather than a theory that is high above so that college students have no real sense of gain.

Regarding the artificial intelligence-related technologies currently used in learning/teaching/work, the statistical results show that computer vision (image understanding, three-dimensional vision, dynamic vision, etc.) The contribution proportion is the highest, indicating that AI technology is used most frequently among practitioners; the second is natural language understanding (AI translation, question-answering system, etc.), of which the proportion of students is the highest, indicating that the student group has the highest contribution to AI. The application of technology is mainly concentrated on natural language understanding; the proportion of teachers' application of the five technologies is relatively balanced. The proportion of practitioners who "have a deep understanding of artificial intelligence technology, but have a partial understanding of its application4.2 Experiment results of complex system modeling.

The results in Fig. 7 show that TLFCM has different functions in the data, and one behaves worse than other models. The one is the best in the data. In the experiment, seven cases in the 17 are better than others.

Fig. 6. Technology acceptance radar chart.



Fig. 7. The relationship between Data_Error and RCGA algebra on different FCM models.

From the experimental results in Fig. 8, as Data_Error increases, the ability of CS-FCM to learn FCM decreases. The change due to lack of information or interference by it and C S-FCM can achieve high accuracy over sparse FCM, which shows that our method is robust to noise in time series. With Data_ With the increase of Error, the dataset used by CS-FCM gradually deviates from the actual situation. This can lead to deviations between the patterns learned by the model and the real-world patterns. During the training process, if Data_ There are many errors, and CS-FCM may fall into a state of overfitting or underfitting. Overfitting means that the model is too complex and has a good fitting effect on the training data, but performs poorly on new data; Underfitting indicates that the model is too simple to capture complex patterns in the data. Generalization ability refers to the model's ability to predict new data. Data_ An increase in Error can lead to a decrease in the generalization ability of CS-FCM, as the knowledge learned during training may not be universally applicable.

Fig. 8. The impact of time series relative data length Nt on Model_Error.

Better than the other four methods, Fig. 9 demonstrates that the Data_Error of CS-FCM stays zero in various cases. For 11 of the 30 examples, LASSOFCM outperforms CS-FCM. CS-FCM works better than LASSOCM in the remaining cases as well. According to Model_Error, CS-FCM performs better than dMAGA and ACORD in 29 of the 30 cases and DandC RCGA and RCGA in every instance. In 26 out of 30 examples, CS-FCM beats LASSOFCM and falls short four times. The fact that the standard deviation of CS-FCM is consistently lower than that of alternative techniques suggests that CS-FCM is stable. CS-FCM beats dMAGA in all but one of the 100-node cases. When it comes to CS-FCM performance, the scenario with S=5 and T=4 is not as good as the scenario with S=1 and T=20. The CS-FCM Model_Error is maintained at a low level concurrently.



Fig. 9. Comparison of CS-FCM and other algorithms in running time.

It can further speed up CS-FCM. Although our algorithm works well in large-scale fuzzy cognitive graph learning problems, several shortcomings still need to be improved. For example, the algorithm in this chapter does not work well on noisy data. The algorithm in this chapter is worth thinking about. After a series of experimental verifications, the CS-FCM algorithm has shown good performance in processing artificial data generated by Gaussian functions. Compared with traditional FCM algorithms, CS-FCM has significantly improved accuracy, boundary clarity, and noise suppression. In addition, CS-FCM is insensitive to the selection of initial parameters and has strong robustness. However, there are also some potential areas for improvement in the CS-FCM algorithm. For example, how to further improve the real-time performance of algorithms to better apply them to real-time image processing systems; How to better handle more complex image data, such as color images or images with more complex textures.

## V. CONCLUSION

The revolution was initiated by artificial intelligence in one area. The proposal of intelligent education points out the direction for development and AI empowerment, and one should also keep pace with the times and use artificial intelligence to empower college students' ideological and political education. It first takes artificial intelligence as the starting point, new forms, models, and paths. Through the questionnaire survey, the value demands empowered by artificial intelligence can be understood. It is possible by artificial intelligence to clarify such persistence and change empowered by artificial intelligence and finally propose an innovative path. However, there are also many potential risks and practical problems at the value premise, technical limits, and specific operation level, such as privacy protection and ideological security risks, the loss of educational subjectivity, the digitization of educational relations, and the lack of specialized talents. Therefore, it is necessary to look at the technical momentum and potential risks of artificial intelligence dialectically, promote the rationality of educational value, strengthen technical supervision, forge an intelligent education team, reasonably define the integration boundary and application scope of artificial intelligence, and combine the main initiative of human beings with the intelligence of machines. It combines strengths, actively explores the path of coexistence and co-prosperity between education and technology, and consciously constructs an intelligent form of them. To model large-scale complex systems, the graph learning method is proposed. The graph from complex system representation data is an urgent problem. Due to the large search space and slow convergence speed, it will have large-scale fuzzy cognitive graphs. Combined with the sparsity of fuzzy cognitive graph, the algorithm decomposes the learning task of fuzzy cognitive graph. The precise recovery ability of compressive sensing for sparse signals allows the algorithm to learn fuzzy cognitive maps accurately. The algorithm breaks through the shortcomings of existing methods and increases the learning scale of FCM from 40 to 1000 nodes, achieving breakthroughs of multiple orders of magnitude.

However, this article also has some limitations. The study of the value demands of ideological and political education for large-scale fuzzy cognition college students requires processing a large amount of data. This includes student personal information, learning behavior, social networks, and other information. Due to the large amount of data, effective data processing and analysis is a huge challenge. The implementation and maintenance of these algorithms and models require a high level of technical proficiency, as well as a significant amount of computing resources and time.

In the future, ideological and political education will be integrated with disciplines such as psychology, sociology, and computer science to study and understand students' thoughts, behaviors, and values from multiple perspectives. This interdisciplinary research method helps to deeply explore the potential of learning the value demands of ideological and political education for large-scale fuzzy cognitive college students.

## COMPETING OF INTERESTS

The authors declare no competing of interests.

## AUTHORSHIP CONTRIBUTION STATEMENT

Chao Xu: Writing-Original draft preparation, Conceptualization, Supervision, Project administration.

Lin Wu: Methodology, Software, Validation.

## DATA AVAILABILITY

On Request

## DECLARATIONS

Not applicable

## REFERENCES

[1] K. A. Walker et al., "Association of peripheral inflammatory markers with connectivity in large-scale functional brain networks of non-demented older adults," Brain Behav Immun, vol. 87, pp. 388–396, 2020.

[2] N. Kampa, R. Scherer, S. Saß, and S. Schipolowski, "The relation between science achievement and general cognitive abilities in large-scale assessments," Intelligence, vol. 86, p. 101529, 2021.

[3] S. Asadianfam, M. Shamsi, and A. R. Kenari, "TVD-MRDL: traffic violation detection system using MapReduce-based deep learning for large-scale data," Multimed Tools Appl, vol. 80, pp. 2489–2516, 2021.

[4] O. Marena, M. S. N. Fitri, O. A. Hisam, A. Kamil, and Z. M. Latif, "Accuracy Assessment of Large-Scale Topographic Feature Extraction Using High Resolution Raster Image and Artificial Intelligence Method," in IOP Conference Series: Earth and Environmental Science, IOP Publishing, 2021, p. 012045.

[5] L. D. Anderson et al., "Large-scale Map of Millimeter-wavelength Hydrogen Radio Recombination Lines around a Young Massive Star Cluster," Astrophys J Lett, vol. 844, no. 2, p. L25, 2017.

[6] M. Lam et al., "Identifying Nootropic Drug Targets via Large-Scale Cognitive GWAS and Transcriptomics," bioRxiv, pp. 2020–2022, 2020.

[7]    O. A. Gansser and C. S. Reich, "A new acceptance model for artificial intelligence with extensions to UTAUT2: An empirical study in three segments of application," Technol Soc, vol. 65, p. 101535, 2021.

[8]    S. A. Simpson and T. S. Cook, "Artificial intelligence and the trainee experience in radiology," Journal of the American College of Radiology, vol. 17, no. 11, pp. 1388–1393, 2020.

[9]    N. Martiniello et al., "Artificial intelligence for students in postsecondary education: a world of opportunity," AI Matters, vol. 6, no. 3, pp. 17–29, 2021.

[10]   P. J. Slanetz, D. Daye, P.-H. Chen, and L. R. Salkowski, "Artificial intelligence and machine learning in radiology education is ready for prime time," Journal of the American College of Radiology, vol. 17, no. 12, pp. 1705–1707, 2020.

[11]   J. Wen, "Innovative application of artificial intelligence technology in college physical education," in Journal of Physics: Conference Series, IOP Publishing, 2021, p. 042028.

[12]   N. Prakash, B. Vaikundaselvan, and S. S. Sivaraju, "Short-Term Load Forcasting for Smart Power Systems Using Swarm Intelligence Algorithm," Journal of Circuits, Systems and Computers, vol. 31, no. 11, p. 2250189, 2022.

[13]   V. Sivalenka, S. Aluvala, N. Fatima, D. R. Kumari, and C. H. Sandeep, "Exploiting Artificial Intelligence to Enhance Healthcare Sector," in IOP Conference Series: Materials Science and Engineering, IOP Publishing, 2020, p. 022061.

[14]   H. Ahmed and L. Devoto, "Undergraduate medical education and the future of surgery," Postgrad Med J, vol. 98, no. e3, pp. e148–e148, 2022.

[15]   M. Tomko, W. Newstetter, M. W. Alemán, R. L. Nagel, and J. Linsey, "Academic makerspaces as a 'design journey': developing a learning model for how women students tap into their 'toolbox of design,'" AI EDAM, vol. 34, no. 3, pp. 363–373, 2020.

[16]   T. Chen, "Research on the dilemma and breakthrough path of ideological and political education in colleges and universities in the era of big data," Journal of Higher Education Research, vol. 3, no. 2, pp. 203–206, 2022.

[17]   L. Song, S. Feng, and T. Zhang, "The Application of Artificial Intelligence in College Ideological and Political Education," Journal of Frontiers in Educational Research, vol. 1, no. 2, pp. 23–26, 2021.

[18]   X. Sun and Y. Zhang, "Research on the framework of university ideological and political education management system based on artificial intelligence," Journal of Intelligent & Fuzzy Systems, no. Preprint, pp. 1–10, 2021.

[19]   A. Kukulska-Hulme, "How should the higher education workforce adapt to advancements in technology for teaching and learning?," Internet High Educ, vol. 15, no. 4, pp. 247–254, 2012.

[20]   U. F. Mustapha, A. Alhassan, D. Jiang, and G. Li, "Sustainable aquaculture development: a review on the roles of cloud computing, internet of things and artificial intelligence (CIA)," Rev Aquac, vol. 13, no. 4, pp. 2076–2091, 2021.

[21]   J. He, S. L. Baxter, J. Xu, J. Xu, X. Zhou, and K. Zhang, "The practical implementation of artificial intelligence technologies in medicine," Nat Med, vol. 25, no. 1, pp. 30–36, 2019.

[22]   R. Dilmurod and A. Fazliddin, "Prospects for the introduction of artificial intelligence technologies in higher education," ACADEMICIA: an international multidisciplinary research journal, vol. 11, no. 2, pp. 929–934, 2021.

[23]   R. Abduljabbar, H. Dia, S. Liyanage, and S. A. Bagloee, "Applications of artificial intelligence in transport: An overview," Sustainability, vol. 11, no. 1, p. 189, 2019.

[24]   D. S. Smys, D. H. Wang, and D. A. Basar, "5G network simulation in smart cities using neural network algorithm," Journal of Artificial Intelligence and capsule networks, vol. 3, no. 1, pp. 43–52, 2021.

[25]   H. H. Tesfamikael, A. Fray, I. Mengsteab, A. Semere, and Z. Amanuel, "Simulation of eye tracking control based electric wheelchair construction by image segmentation algorithm," Journal of Innovative Image Processing (JIIP), vol. 3, no. 01, pp. 21–35, 2021.

[26]   Z. B. Jimenez et al., "Matrix representation and simulation algorithm of spiking neural P systems with structural plasticity," Journal of Membrane Computing, vol. 1, pp. 145–160, 2019.

[27]   S. Harchaoui and P. Chatzimpiros, "Energy, Nitrogen, and Farm Surplus Transitions in Agriculture from Historical Data Modeling. France, 1882–2013.," J Ind Ecol, vol. 23, no. 2, pp. 412–425, 2019.

[28]   O. Melnychenko, "Application of artificial intelligence in control systems of economic activity," Virtual Economics, vol. 2, no. 3, pp. 30–40, 2019.

[29]   J. Chen, K. Li, Z. Zhang, K. Li, and P. S. Yu, "A survey on applications of artificial intelligence in fighting against COVID-19," ACM Computing Surveys (CSUR), vol. 54, no. 8, pp. 1–32, 2021.

[30]   K. Zarina I, B. Ildar R, and S. Elina L, "Artificial Intelligence and Problems of Ensuring Cyber Security.," International Journal of Cyber Criminology, vol. 13, no. 2, 2019.

# Predicting Students' Academic Performance Through Machine Learning Classifiers: A Study Employing the Naive Bayes Classifier (NBC)

Xin ZHENG[1], Conghui LI[2]

College of Artificial Intelligence, North China University of Science and Technology, Tangshan 063210, China[1]
College of Management and Economics, North China University of Science and Technology, Tangshan 063210, China[2]

*Abstract*—**Modern universities must strategically analyze and manage student performance, utilizing knowledge discovery and data mining to extract valuable insights and enhance efficiency. Educational Data Mining (EDM) is a theory-oriented approach in academic settings that integrates computational methods to improve academic performance and faculty management. Machine learning algorithms are essential for knowledge discovery, enabling accurate performance prediction and early student identification, with classification being a widely applied method in predicting student performance based on various traits. Utilizing the Naive Bayes classifier (NBC) model, this research predicts student performance by harnessing the robust capabilities inherent in this classification tool. To bolster both efficiency and accuracy, the model integrates two optimization algorithms, namely Jellyfish Search Optimizer (JSO) and Artificial Rabbits Optimization (ARO). This underscores the research's commitment to employing cutting-edge machine learning and algorithms inspired by nature to achieve heightened precision in predicting student performance through the refinement of decision-making and prediction quality. To classify and predict G1 and G3 grades and evaluate students' performance in this study, a comprehensive analysis of the information pertaining to 395 students has been conducted. The results indicate that in predicting G1, the NBAR model, with an F1_Score of 0.882, performed almost 1.03% better than the NBJS model, which had an F1_Score of 0.873. In G3 prediction, the NBAR model outperformed the NBJS model with F1_Score values of 0.893 and 0.884, respectively.**

*Keywords*—*Machine learning; Naive Bayes Classifier; Artificial Rabbits Optimization; Jellyfish Search Optimizer; student performance*

## I. INTRODUCTION

Education is the foundational pillar for any nation or society, embodying a crucial element that provides guidance, societal status, extensive knowledge, and avenues for exploration [1–3]. Modern universities must analyze performance, identify uniqueness, and develop a strategic plan. Management should prioritize understanding admitted students' diverse characteristics. In the competitive academic landscape, excellence in student performance is crucial for higher learning institutions [4–6]. Knowledge discovery (KD) involves extracting meaningful, unknown, and potentially valuable information from extensive databases. Data mining ($DM$) is crucial for educational data analysis, offering various methods for this purpose [7, 8]. The substantial amount of student data

in databases exceeds the human capacity for manual analysis, necessitating automated techniques [9]. This includes creating early warning systems to reduce costs, save time, and optimize resources [10, 11]. Educational Data Mining ($EDM$) is a theory-oriented $DM$ approach employed in academic and educational settings. Its primary goal is to create computational methods that integrate theory and data. The objective of EDM is to improve and support academic performance among students and graduates and to enhance the management of faculty information within educational institutions [12–14].

Machine learning (ML) algorithms serve as indispensable tools for knowledge discovery, playing a pivotal role in various applications. One of their crucial functions lies in accurate performance prediction, a capability that proves instrumental in identifying struggling students early. By leveraging these algorithms, educational institutions can proactively address academic challenges, fostering a more supportive and responsive learning environment [15]. Various machine learning methods, including classification [16], prediction [17, 18], and clustering [19], are continuously evolving and expanding the scope of data mining.

Classification, the most common and effective data mining approach for categorizing and predicting values, also applies to EDM [20]. In the context of student performance prediction, classification refers to grouping or classifying pupils based on specific traits or features. These traits include past academic achievement, demographic information, socioeconomic background, study habits, and other pertinent data. Patterns and links within the data are detected using classification algorithms, allowing predictions about future student performance to be generated [21–24]. For instance, using the ICRM classifier, Marquez-Vera et al. [25] tackled the intricate task of predicting student failure in academic contexts by utilizing a genetic programming algorithm and diverse $DM$ approaches. Employing real data from 670 high school students in Zacatecas, Mexico, the research addressed challenges such as high dimensionality and imbalanced datasets. It strategically selected influential attributes, rebalances data, and incorporated cost-sensitive classification methods. Additionally, the study introduced a genetic programming model, comparing its interpretability and accuracy with other white box techniques. The ultimate goal was to identify the most effective approach for enhancing classification accuracy, particularly in predicting students at

risk of failure. The findings contributed valuable insights to predictive modeling in education, offering nuanced perspectives on factors influencing student outcomes and guiding targeted interventions for improved educational support. Hu and Song [26] focused on the analysis and evaluation of student achievement as a crucial aspect of teaching and school management. The objective was to scientifically assess academic performance, providing accurate insights for teachers and enabling students to understand their learning situation. The research utilized the XGBoost algorithm for classifying and evaluating student performance through statistical analysis of basic data. A performance evaluation model was established, taking into account curriculum relevance by statistically compiling student performance data. The subjective and objective structural entropy weight method was employed to classify characteristic importance results, offering insights into relevant courses. Moreover, the XGBoost method was used to predict grades for unfinished courses based on completed course results. The study aimed to comprehensively, objectively, and reasonably evaluate students' learning situations, contributing valuable information for teaching management and improvement strategies. Kabakchieva et al. [27] aimed the outcomes of a data mining research conducted at a prestigious Bulgarian university. The primary objective was to showcase the substantial potential of data mining applications in university management, particularly in optimizing enrollment campaigns and attracting high-caliber students. The research focused on developing data mining models to predict student performance, utilizing personal, pre-university, and university-performance characteristics. The dataset encompassed information about students admitted to the university over three consecutive years. Various well-known data mining classification algorithms, including a rule learner, a decision tree classifier, a neural network, and a Nearest Neighbour classifier, were applied and their performances were analyzed and compared. The study aimed to contribute insights for more effective university management by employing data mining techniques to predict and understand student performance. By presenting a model for predicting poor academic performance among first-year students, Tamasiri et al. [28] aimed to address the challenging task of predicting student attrition, particularly dealing with class-imbalanced data common in the realm of student retention. The study contrasted four widely used classification techniques logistic regression, decision trees, neural networks, and support vector machines with three alternative data balancing strategies: over-sampling, under-sampling, and synthetic minority over-sampling ($SMOTE$). The research aimed to retain overall excellent classification performance while improving predicting accuracy for the minority class using large-scale institutional student data from 2005 to 2011. Based on the 10-fold holdout sample, the support vector machine and SMOTE data-balancing approach produced the best classification result, with an overall accuracy of 90.24% for the minority class, the three data-balancing strategies increased prediction accuracy. Additionally, sensitivity analyses identified crucial variables for accurately predicting student attrition, suggesting the potential application of these models to reduce student dropout rates by accurately identifying at-risk students. Marbouti et al. [29] utilized

predictive modeling methods to early identify at-risk students in courses employing standards-based grading. The goal of the study was to modify at-risk prediction models to take use of standards-based grading, which offers advantages over traditional score-based grading in education. Prediction approaches were limited to using the course instructors' access to performance data from the previous semester. The study prioritized minimizing false negatives (type II errors) in identifying at-risk students without significantly increasing false positives (type I errors). To enhance generalizability and accuracy, a feature selection method was applied to reduce the number of variables in each model. Among the seven tested modeling methods, the Naive Bayes ($NB$) Classifier model and an Ensemble model showed the most promising results, contributing insights for more effective educational interventions in identifying at-risk students.

While various classification algorithms have received considerable attention in recent studies, Naive Bayes Classifier (NBC) has been comparatively less explored. This research introduces and evaluates NBC alongside two hybrid models optimized using Jellyfish Search Optimizer (JSO) and Artificial Rabbits Optimization (ARO). The study comprehensively assesses their estimation capabilities by training 70% of the models on literature-derived input parameters and testing the remaining 30%, enabling comparisons with other models and evaluations of enhanced versions of a single model. The examination involves statistical metrics in two distinct phases and categorizing students into four grade classes, providing a thorough comparative analysis. Ultimately, the study identifies the optimal model for understanding and anticipating students' academic performance, emphasizing NBC's adaptability, uncertainty estimation, and interpretability when integrated with optimization algorithms to enhance predictive accuracy and improve educational outcomes. In the following Section II outlines the methods, procedures, and details of your research, encompassing data collection, experimental design, participants, materials, and any statistical or computational methods employed. Moreover, the Experimental Design or Data Collection subsection provides a detailed account, including variables, controls, and procedures implemented. Dataset overview is given in Section III. Section IV presents study findings using tables, graphs, or figures, incorporating both descriptive and inferential statistical analyses. In Section V, results are interpreted in the context of addressing implications, limitations, and potential future research directions. Finally Section VI summarizes main findings, emphasizing their significance, broader implications, and suggesting potential applications or areas for further investigation.

## II.    METHODOLOGY AND STRATEGY

### A.  Naive Bayes Classifier (NBC)

The NBC is a probabilistic classifier that utilizes Bayes' theorem under the assumption of high independence. This algorithm was formulated by Thomas Bayes, a British scientist. The theoretical foundation of NBC revolves around predicting future opportunities by drawing on past experiences [30]. Assume a category of documents $D = \{d_1, d_2, ..., d_n\}$ and $m$

potential classes $C = \{c_1, c_2, \ldots, c_m\}$. Where $W = \{w_1, w_2, \ldots, w_s\}$ represent the collection of distinct terms present in at least one of the documents in $D$. The subsequent formula can be used to calculate the probability that a document $d$ belongs to class $c$ [31].

$$P(c|d) = \frac{P(d|c)\, P(c)}{P(d)} \qquad (1)$$

The denominator of Eq. (1) is often omitted in Maximum *A* Posteriori ($MAP$) calculation for parametric numerical problems because $P(d)$ is a constant for the known data set size. In the context of a $NBs$ model, it accepts that each term or word, $w_k$, independently occurs in a document given the class $c$. Consequently, Eq. (1) is simplified to reflect this assumption:

$$P(c|d) \propto P(c) \prod_{k=1}^{n_d} \left[P(w_k \,|\, c)\right]^{t_k} \qquad (2)$$

In the provided context, $n_d$ denotes the count of single words in file $d$, and $t_k$ represents the frequency of each word $w_k$. To address concerns regarding floating$-$point underflow, an alternative equation is utilized:

$$\log P(c|d) \propto \log P(c) + \sum_{k=1}^{n_d} \left[t_k \log P(w_k \,|\, c)\right] \qquad (3)$$

The classification of document $d$ is determined as the class $c^*$ that maximizes the logarithm of $P(c|d)$ in Eq. (3).

$$c^* = argmax_c \in \mathbb{C} \left\{ \log P(c) + \sum_{k=1}^{n_d} \left[t_k \log P(w_k \,|\, c)\right] \right\} \qquad (4)$$

In the application of the Naive Bayes classifier (NBC), $P(c)$ and $P(w_k \,|\, c)$ can derive estimations as follows:

$$\hat{P}(c) = \frac{N_c}{N} \; and \; \hat{P}(w_k|c) = \frac{N_{w_k}}{\sum_{w_i \in \mathbb{W}} N_{w_i}} \qquad (5)$$

Where $N$ signifies the total document count, $N_c$ denotes the quantity of records in class $c$, and $N_{w_i}$ represents the word's frequency $w_i$ in class $c$. Utilizing these approximations, the computation of the expression on the right$-$hand side of Eq. (4) essentially becomes a counting challenge.

*B. Jellyfish Search Optimizer ($JSO$)*

One of the contemporary swarm-based metaheuristics is the JSO, introduced by Chou and Truong in 2021 [32]. The $JSO$ algorithm emulates the foraging behaviour of jellyfish in search of oceanic sustenance [33].

*1) Mathematical model:* The JSO algorithm adheres to three ideal principles:

*a) Marine flow:* To locate and feed on smaller planktonic organisms, jellyfish utilize Eq. (6) to detect the direction of ocean currents.

$$\vec{O} = X' - \beta \times M \times r(0,1) \qquad (6)$$

The direction of the ocean current is denoted as $\vec{O}$, where $\beta$ $(\beta > 0)$ represents the length distribution coefficient of $\vec{O}$.

The current best location of the jellyfish swarm is denoted as $X'$, and $M$ signifies the mean location of all jellyfish.

The new location of each jellyfish can be defined as follows:

$$X_i(t + 1) = X_i(t) + r(0,1) \times \vec{O} \qquad (7)$$

Following the adjustment of each jellyfish's position, the current location of the jellyfish is chosen as a preferred destination, potentially representing a position with increased accessibility to food sources.

*b) Blooming of jellyfish:* Jellyfish within a bloom display two types of motion: passive and active. The subsequent section introduces the mathematical models that characterize these waves:

Passive waves: 
$$X_i(t + 1) = X_i(t) + \lambda \times r(0,1) \times (w_b - L_b) \qquad (8)$$

$\lambda$ $(\lambda > 0)$ signifies a coefficient associated with the extent of passive waves. $w_b$ and $L_b$ denote the lower and upper bounds of the search space, correspondingly.

Active waves: 
$$X_i(t + 1) = X_i(t) + r(0,1) \times \vec{D} \qquad (9)$$

which

$$\vec{D} = \begin{cases} X_i(t) - X_j(t) & if \quad g(X_i) < \; g(X_j) \\ X_j(t) - X_i(t) & if \quad g(X_i) \geq \; g(X_j) \end{cases} \qquad (10)$$

The functions $g(X_i)$ and $g(X_j)$ represent the objective purpose values corresponding to jellyfish $i$ and $j$, correspondingly.

*c) Temporal regulation mechanism:* Within this framework, a mechanism for time control is utilized to regulate both the motion of jellyfish within the bloom and their navigation toward ocean currents. The temporal control function is denoted as:

$$T(t) = \left| \left(1 - \frac{t}{MaxIter}\right) \times (2 \times r(0,1) - 1) \right| \qquad (11)$$

where, $t$ shows the time index given as the iteration quantity and $MaxIter$ represents the iteration $max$ number. $t$ denotes the time index, representing the iteration number, and $MaxIter$ signifies the maximum number of iterations.

*2) Population initialization:* In this optimizer, the initial population is generated using the Logistic map.

$$X_{i+1} = \upsilon X_i(1 - X_i), \qquad 0 \leq X_0 \leq 1 \qquad (12)$$

$X_i$ and $X_0$ represent the positional values for the $i_{th}$ jellyfish and a randomly selected place, respectively. It is set to 4 in all testing.

*a) Boundary handling mechanism:* If a jellyfish surpasses the confines of the specified search space, it will be realigned within those limits, according to Eq. (13).

$$\begin{cases} X'_{i,d} = (X_{i,d} - W_{b,d}) + L_{b,d} & if \quad X_{i,d} > \; W_{b,d} \\ X'_{i,d} = (X_{i,d} - L_{b,d}) + W_{b,d} & if \quad X_{i,d} < \; W_{b,d} \end{cases} \qquad (13)$$

$X_{i,d}$ and $X'_{i,d}$ represent the current and updated location of the $d_{th}$ dimension of the $i_{th}$ jellyfish, respectively. $W_{b,d}$ and $L_{b,d}$ symbolizes the lower and upper bounds of the $d_{th}$ dimension within the search space.

The flowchart illustrating the JSO process is shown in Fig. 1.



Fig. 1. Flowchart of JSO.

## C. Artificial Rabbits Optimization (ARO)

Rabbits' adoption of survival strategies within their natural environment has significantly influenced the formulation of ARO. These strategies are designed to effectively counteract predators and optimize the rabbit's ability to evade surveillance. ARO, in its design, integrates the rabbit's dual strategies of hunting and hiding, along with its adept energy management techniques, to seamlessly transition between these adaptive approaches [34].

*1) Detour foraging:* Rabbits commonly employ a circuitous method during their quest for food, emphasizing distant food sources while overlooking nearby ones. Consider an ARO scenario in which each rabbit in a group has its

territory complete with caves and grass. Serendipitous encounters with each other's feeding areas are common. A mathematical framework is presented that captures rabbits' departure search actions:

$$\vec{B}_i(t+1) = x_j(t) + S \times (x_i(t) - x_j(t)) + w(0.5 \times (0.05 + r_1)) \times m_1, \quad (14)$$

$$i, j = 1, \dots, n \text{ and } j \neq 1$$

$$S = M \times v \quad (15)$$

$$M = (e - e^{\left(\frac{t-1}{I}\right)^2}) \times \sin(2\pi r_2) \quad (16)$$

$$v(y) = \begin{cases} 1 & if \quad y = f(1) \\ 0 & else \end{cases} \quad k = 1, \dots, d \text{ and } l = 1, \dots, \lceil r_3 \times d \rceil \quad (17)$$

$$f = p(d) \quad (18)$$

$$m_1 = N(0,1) \quad (19)$$

The population size of rabbits is denoted as *n*, and the dimensions of the problem are denoted as *d*.

$\vec{B}_i(t+1)$ signifies the standard normal distribution describes the distribution of the $i\_th$ rabbit's location at times $t+1$ and $n1$. $T$ represents the $max$ number of iterations. $x_i(t)$ signifies the place of the $i_{th}$ rabbit at time $t$. A random permutation of numbers between 1 and $d$ is produced by the function $p$.

$w$ functions as a tool inside the algorithm, promoting the diverse collection of components from the traveler to introduce variety into the process. $r_1$, $r_2$, and $r_3$ depict random numbers within the (0,1) range. The variable $s$ denotes the run length, indicating the pace of development during reroute scavenging.

*2) Random hiding:* Rabbits tend to randomly select a burrow to seek shelter, a behaviour crucial for their survival. The mathematical model elucidating this performance is articulated through the equations presented below. The formulation for the $j_{th}$ burrow of the $i_{th}$ rabbit is expressed as:

$$\vec{B}_i(t+1) = x_i(t) + N \times f \times \vec{x}_i(t), \quad i, j = 1, \dots, n \text{ and } j \neq 1 \quad (20)$$

$$D = \frac{I - t + 1}{I} \times r_4 \quad (21)$$

$$m_2 = N(0,1) \quad (22)$$

$$f(y) = \begin{cases} 1 & if \quad y = g(1) \\ 0 & else \end{cases} \quad k = 1, \dots, d \quad (23)$$

$$\vec{R}_{i,r}(t) = \vec{x}_i(t) + N \times f \times \vec{x}_i(t) \quad (24)$$

N embodies the concealment parameter, undergoing a linear reduction from 1 to $\frac{1}{I}$ throughout an iterative process that incorporates random perturbations.

In the eventual scenario where either detour foraging strategies or random hiding methods are employed, the update of the $i_{th}$ rabbit's location observes to the Eq. (25):

$$\vec{x}_i(t+1)$$

$$= \begin{cases} \vec{x}_i(t) & g(\vec{x}_i(t)) \leq g\left(\vec{B}_i(t+1)\right) \\ \vec{B}_i(t+1) & g(\vec{x}_i(t)) > g\left(\vec{B}_i(t+1)\right) \end{cases} \qquad (25)$$

*3) Energy shrinks:* While the rabbit persists in its cyclic behaviour of detouring to find food and intermittently hiding at random, its energy level regularly diminishes. Hence, the integration of energy factors becomes critical within the ARO framework:

$$E(t) = 4\left(1 - \frac{t}{l}\right) ln\frac{1}{r} \qquad (26)$$

Fig. 2 illustrates the flowchart of ARO, and Algorithm 1 provides its pseudo-code.



Fig. 2. Flowchart of ARO.

| Algorithm 1: Pseudo-Code of *ARO* Algorithm |
|---|
| Randomly initialize a set of rabbits. $X_i$ (solutions) and evaluate their fitness <br> While the stop criterion is not satisfied, do <br> for each individual $X_i$ do <br> Compute the energy factor A using Eq. (26). <br> if A > 1 <br> Choose a rabbit randomly from other individuals. <br> Compute R using Eqs. (15) − (19). <br> Perform detour foraging using Eq. (14). <br> Compute the fitness $Fit_i$. <br> Upgrade the position of the current individual using Eq. (25). <br> else <br> Generate d burrows and randomly pick one as hiding using Eq. (24). <br> Perform random hiding using Eq. (20). <br> Compute the fitness $Fit_i$. <br> Update the position of the individual using Eq. (25). <br> end if <br> Upgrade the best solution found so far $X_{best}$ <br> end for <br> end while <br> return $X_{best}$ |

## III. DATASET OVERVIEW

### A. Data Preparation

Data mining is a strategic business process, systematically delving into vast datasets to unearth significant patterns and rules that contribute valuable insights [35]. Classification and regression are pivotal objectives in data mining, playing crucial roles in extracting valuable insights from complex datasets by identifying meaningful patterns and relationships [36]. The primary distinction lies in the output representation, with classification producing discrete results and regression generating continuous outcomes. Evaluation metrics also differ, as classification models are often assessed using the percentage of correct classifications; while regression commonly employs the root mean squared metric [37].

This study aims to develop a robust method for accurately assessing students' academic performance and contextual factors. The dataset undergoes essential preprocessing, including transforming text into numerical values. Attributes are selected to describe performance based on individual information and academic conditions, utilizing two questionnaire methods and academic histories. The dataset encompasses a varied array of variables that have the potential to influence students' academic outcomes. The dataset incorporates information, including students' school, gender, age, home address, family size, parental cohabitation status, as well as details about the education and occupations of both parents. The dataset also includes information on factors influencing school choice, such as proximity, reputation, course preference, and others. It covers details about the student's guardian, weekly study time, travel time to school, past class failures, participation in educational support, family educational support, involvement in paid classes and extracurricular activities, internet access at home, aspirations for higher education, nursery school attendance, engagement in romantic relationships, family relationship quality, current health status, socializing with friends, weekday and weekend

alcohol consumption, free time after school, and the number of school absences. These diverse input variables, encompassing nominal, numeric, and binary data types, provide a comprehensive and varied source of information for the study. In addition to these inherent traits, three supplementary variables, namely, $G2$, $G1$, and $G3$, depict students' grades across three assessment periods throughout their academic journey, spanning from zero (indicating the lowest grade) to $20$ (representing the highest grade attainable). To classify their scores, a segmentation was applied, dividing them into four categories: 0–12 denoting Poor performance, 12–14 indicating Acceptable, 14–16 representing Good, and 16–20 reserved for Excellent academic achievements.

In this research, a correlation matrix encompassing the examined input and output variables is depicted in Fig. 3. Parents' educational background, particularly the mother, exerted the most positive influence on students' scores in G1 and G3, with the father's education also demonstrating effectiveness. As anticipated, family support and aspirations for higher education had positive effects on outcomes, while the influence of prior student failures was negative. Additionally, there was a noticeable gender effect on scores in G1 and G3.

### B. Evaluation of Models' Applicability

In classification problems, Accuracy is a widely used metric that gauges overall model performance based on True Positives ($TP$), True Negatives ($TN$), False Positives ($FP$), and False Negatives ($FN$). While Accuracy is common, it has limitations in imbalanced data situations, favoring the majority class and providing limited insights. Three additional metrics, Recall, Precision, and F1-Score, address this. Recall evaluates a model's ability to correctly identify all relevant instances within a specific class, which is crucial for reducing False Negatives. Precision assesses the accuracy of positive predictions, reducing False Positives. The F1-Score combines Precision and Recall, offering a balanced assessment, especially valuable in imbalanced data scenarios.

These metrics, outlined through mathematical formulas Eq. (27) to Eq. (30), collaboratively contribute to a more comprehensive grasp of a classification model's effectiveness [38]. Particularly valuable in tackling class imbalances, they empower researchers and data analysts to make well-informed decisions and adjustments, enhancing model performance in challenging scenarios involving imbalanced data.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{27}$$

$$Precision = \frac{TP}{TP + FP} \tag{28}$$

$$Recall = TPR = \frac{TP}{P} = \frac{TP}{TP + FN} \tag{29}$$

$$F1\_score = \frac{2 \times Recall \times Precision}{Recall + Precision} \tag{30}$$

Fig. 3. Correlation matrix for the input and output variables.

## IV. RESULTS

To improve the accuracy of the NBC model in predicting G1 and G3, this study employed two optimization algorithms, ARO and JSO. The dataset was divided, with 70% allocated for the training phase and the remaining 30% for thorough testing, enabling a comprehensive assessment of predictive capabilities. The data underwent processing after a detailed evaluation of the models' classification during training and testing, involving 395 students and grounded in their test results (specifically G1 and G3 values).

The primary objective involved fine-tuning and optimizing model parameters through these algorithms. To assess the convergence of these optimization methods, a convergence curve, depicted in Fig. 4, tracked accuracy over 200 iterations. The convergence rate of the NBAR model in G1 and G3 is similar, with a noticeable shift to a linear pattern around the 150th iteration in the convergence process. In contrast, the NBJS model, which predicts both G1 and G3 metrics, is the optimal model, achieving superior accuracy levels before the 150th iteration.



Fig. 4. Convergence of hybrid models based on ribbon plot.

This section assesses how each model contributes to predicting students' academic performance based on the G1 and G3 grades. Table I presents Accuracy, Recall, Precision, and F1-score measures for the training and testing phases across all models. Notably, the NBAR model demonstrated superior performance, exhibiting higher metric values. Specifically, in both G1 and G3, the NBAR model achieved maximum metric values with Accuracy=0.889, Precision = 0.885, Recall = 0.886, and F1-score = 0.882 for G1, and Accuracy= 0.894, Precision= 0.893, Recall= 0.894, and F1-score= 0.893 for G3, respectively. Additionally, the waterfall plot in Fig. 5 provides a visual assessment of the performance of the presented models.

The students were classified into four distinct groups based on their scores: Poor (0 to 12), Acceptable (12 to 14), Good (14 to 16), and Excellent (16 to 20). In terms of Precision within G1 estimation, shown in Table II, the NBJS model showcases superior performance, achieving values of 0.923 and 0.769 in the Excellent and Good categories, respectively. Conversely, for the Acceptable and Poor groups, the NBAR model demonstrates Precision values of 0.880 and 0.912, respectively. As indicated in Table III for G3 precision values, the NBAR model demonstrates superior performance in the Excellent and Acceptable groups with values of 0.917 and 0.768, respectively. Conversely, the NBJS model, with values of 0.836 and 0.957, is more suitable for the Good and Poor groups.

TABLE I.        RESULT OF PRESENTED MODELS

| | Model | phase | Index values | | | |
|---|---|---|---|---|---|---|
| | | | Accuracy | Precision | Recall | F1 _Score |
| G1 | NBC | Train | 0.884 | 0.882 | 0.885 | 0.881 |
| | | Test | 0.822 | 0.810 | 0.822 | 0.811 |
| | | All | 0.866 | 0.861 | 0.866 | 0.861 |
| | NBAR | Train | 0.917 | 0.921 | 0.917 | 0.915 |
| | | Test | 0.814 | 0.801 | 0.814 | 0.801 |
| | | All | 0.889 | 0.885 | 0.886 | 0.882 |
| | NBJS | Train | 0.899 | 0.900 | 0.899 | 0.894 |
| | | Test | 0.831 | 0.825 | 0.831 | 0.821 |
| | | All | 0.878 | 0.877 | 0.879 | 0.873 |
| G3 | NBC | Train | 0.866 | 0.866 | 0.866 | 0.865 |
| | | Test | 0.898 | 0.905 | 0.898 | 0.900 |
| | | All | 0.876 | 0.878 | 0.876 | 0.876 |
| | NBAR | Train | 0.883 | 0.882 | 0.883 | 0.882 |
| | | Test | 0.915 | 0.918 | 0.915 | 0.916 |
| | | All | 0.894 | 0.893 | 0.894 | 0.893 |
| | NBJS | Train | 0.870 | 0.869 | 0.870 | 0.869 |
| | | Test | 0.915 | 0.922 | 0.915 | 0.917 |
| | | All | 0.884 | 0.885 | 0.884 | 0.884 |



NBAR-G1



NBAR-G3

Fig. 5. Waterfall plot utilized to assess the performance of the presented models.

TABLE II. ASSESSMENT METRICS FOR THE PERFORMANCE OF THE GENERATED MODELS DERIVED FROM G1 PREDICTION

| Model | Grade | Index values | | |
|-------|-------|-----------|--------|----------|
| | | *Precision* | *Recall* | *F1-score* |
| NBC | Excellent | 0.875 | 0.854 | 0.864 |
| | Good | 0.745 | 0.704 | 0.724 |
| | Acceptable | 0.800 | 0.647 | 0.715 |
| | Poor | 0.904 | 0.970 | 0.936 |
| NBAR | Excellent | 0.897 | 0.854 | 0.875 |
| | Good | 0.768 | 0.796 | 0.782 |
| | Acceptable | 0.880 | 0.647 | 0.746 |
| | Poor | 0.912 | 0.983 | 0.946 |
| NBJS | Excellent | 0.923 | 0.878 | 0.900 |
| | Good | 0.769 | 0.741 | 0.755 |
| | Acceptable | 0.875 | 0.617 | 0.724 |
| | Poor | 0.895 | 0.987 | 0.939 |

TABLE III.    ASSESSMENT METRICS FOR THE PERFORMANCE OF THE GENERATED MODELS DERIVED FROM G3 PREDICTION

| Model | Grade | Index values | | |
|---|---|---|---|---|
| | | *Precision* | *Recall* | *F1-score* |
| NBC | Excellent | 0.909 | 0.750 | 0.822 |
| | Good | 0.772 | 0.733 | 0.752 |
| | Acceptable | 0.691 | 0.758 | 0.723 |
| | Poor | 0.949 | 0.966 | 0.957 |
| NBAR | Excellent | 0.917 | 0.825 | 0.868 |
| | Good | 0.788 | 0.867 | 0.825 |
| | Acceptable | 0.768 | 0.694 | 0.729 |
| | Poor | 0.949 | 0.966 | 0.957 |
| NBJS | Excellent | 0.833 | 0.750 | 0.790 |
| | Good | 0.836 | 0.767 | 0.800 |
| | Acceptable | 0.761 | 0.871 | 0.812 |
| | Poor | 0.957 | 0.957 | 0.957 |

In Fig. 6, 3D walls illustrate a detailed comparison between predicted and measured values, presenting the distribution of students across categories for G1 and G3. Individual graphs for each category (Poor, Acceptable, Good, and Excellent) are included. It is noteworthy that the institute's report mentions a total of 395 students. The subsequent sections undertake a comprehensive assessment of the models' classification effectiveness.

As per the chart, the data for G1 reveals that 232 individuals fall into the Poor category, 68 individuals in the Acceptable category, 54 individuals in the Good category, and 41 individuals in the Excellent category. Notably, the NBJS model is the most effective classifier for the Poor and Excellent segments, demonstrating accurate predictions. Conversely, in the Acceptable and Good groups, the NBAR model outperforms, exhibiting superior performance.

As depicted in the G3 graph, recorded figures for the Poor, Acceptable, Good, and Excellent categories were 233, 62, 60, and 40 students, respectively. An exception arises as the single NBC model, and NBAR exhibit superior performance in the Poor class. Moreover, the NBAR model consistently demonstrates superior performance in the Excellent and Good classes. However, the NBJS model outperforms others in the Acceptable class, while the NBAR model exhibits a noticeable performance drop.

Valuable information regarding the precise classification of students and instances of misclassifications can be extracted from the confusion matrix illustrated in Fig. 7. In G1 estimation, the NBAR model accurately classified a total of 350 students, encompassing 35 Excellent, 43 Good, 44 Acceptable, and 228 Poor students, in their respective grades, while 45 students were misclassified. In contrast, the NBJS model had 48 misclassifications, indicating a 6.25% difference in G1 between the two hybrid models, with the NBAR model exhibiting superior performance. In G3, the NBAR model achieved precise categorization for 353 students, accurately placing them in their respective grades, with only 42 students being misclassified. Similarly, the NBJS model also achieved 353 correct predictions.

Fig. 6.   3D walls for the comparison between predicted and measured values.

Fig. 7. Confusion matrix illustrating the accuracy of the model.

Fig. 8. Result of ROC curve.

The Receiver Operating Characteristics ($ROC$) curve is crucial for assessing classification algorithms. This curve evaluates the model's performance by graphing $TP$ rates against $FP$ rates. A test exhibiting perfect discrimination would manifest as an ROC plot traverse the upper left corner, indicating both 100% sensitivity and 100% specificity. The ROC curve analysis in Fig. 8 shows that, the $NBC$ emerges as the best overall classifier in G1 prediction, particularly for the poor and excellent classes, as evidenced by its proximity to 1 on the ROC curve. In the framework of G3, there is no discernible trend for comparing the performance of the models; however, there is a relative improvement in predicting the Poor group.

## V. DISCUSSION

The study has limitations that should be considered. Firstly, its focus on a specific cohort of 395 students may restrict the generalizability of the findings to a broader student population. Additionally, the evaluation predominantly revolves around quantitative metrics, such as Accuracy, Recall, Precision, and F1-score, potentially overlooking qualitative nuances in student academic performance. The exclusive emphasis on specific grading criteria (G1 and G3) raises questions about the adaptability of the proposed methodology to different grading systems or academic contexts. Despite these limitations, the study represents a notable advancement in predictive modeling for education.

In addition, Future studies in the realm of predictive modeling in education could explore diverse paths to advance the field. Key areas for investigation include assessing the generalizability of hybrid models across different educational settings and student populations, incorporating qualitative factors to provide a more comprehensive understanding of academic performance, and evaluating the adaptability of the proposed methodology to various grading systems. Longitudinal analyses and the integration of real-time data offer opportunities for dynamic predictions and a deeper exploration of academic trajectories. Comparative studies with other predictive models, ethical considerations, and impact assessments on implementation in educational institutions are also important avenues for further research. Addressing these aspects will contribute to refining predictive models and enhancing their practical application in education.

In addition, Table IV is shown to compare the accuracy of the developed best model with other published papers.

TABLE IV.    COMPARISON BETWEEN THE PRESENTED AND PUBLISHED PAPERS

| No. | Paper | Model | Accuracy |
|---|---|---|---|
| 1 | Al-Radaideh et al. [39] | DTC | 87.9% |
| 2 | Bichkar and R. R. Kabra [40] | DTC | 69.94% |
| 3 | Carlos et al. [41] | ADTree | 97.3% |
| 4 | Kabakchieva [42] | DTC | 72.74% |
| 5 | Nguyen and Peter [43] | DTC | 82% |
| 6 | Edin Osmanbegovic et al. [44] | NBC | 76.65% |
| 7 | Present study | NBAR | 89.4% |

## VI. CONCLUSION

This study underscores the vital role of data-driven predictive models in education, emphasizing the need to consider qualitative and quantitative factors in forecasting and assessing student academic performance. The research introduces innovative hybrid models that integrate the Naive Bayes classifier (NBC) with optimization algorithms, namely Jellyfish Search Optimizer (JSO) and Artificial Rabbits Optimization (ARO). The study showcases a cutting-edge methodology demonstrating how the precision and effectiveness of predictive models can be enhanced through advanced machine learning and optimization algorithms. The thorough assessment using essential metrics such as Accuracy,

Precision, Recall, and F1-score underscores these meta-heuristic algorithms' capability to optimize classification results. This study focuses on classifying grades G1 and G3 for a cohort of 395 students. In predicting G1, the NBAR model demonstrated superior performance compared to the NBJS model based on the F1-score criterion, surpassing it by approximately 1.03% and outperforming the NBC model by 2.39%. Regarding Recall, this advantage amounted to 0.18% and 2.83%, respectively. Furthermore, in the G3 forecast, the NBAR model exhibited better performance in the F1-score criterion, approximately 1.01% better than the NBJS model and 1.1% better than the NBC model. This superiority in Recall is represented by percentages of 1.01% and 2.02%, respectively. This research represents a significant advancement in predictive modelling within the field of education, presenting promising avenues for improving the precision and efficiency of evaluating academic performance.

## REFERENCES

[1] T. PanduRanga Vital, B.G. Lakshmi, H. Swapna Rekha, M. DhanaLakshmi, Student Performance Analysis with Using Statistical and Cluster Studies, in: Soft Computing in Data Analytics: Proceedings of International Conference on SCDA 2018, Springer, 2019: pp. 743–757.

[2] A. Behr, M. Giese, H.D. Teguim Kamdjou, K. Theune, Motives for dropping out from higher education—An analysis of bachelor's degree students in Germany, Eur J Educ 56 (2021) 325–343.

[3] D. Buenaño-Fernandez, W. Villegas-CH, S. Luján-Mora, The use of tools of data mining to decision making in engineering education—A systematic mapping study, Computer Applications in Engineering Education 27 (2019) 744–758.

[4] B.A. Pereira, A. Pai, C. Fernandes, A comparative analysis of decision tree algorithms for predicting student's performance, Int J Eng Sci 7 (2017) 10489–10492.

[5] A. Shanthini, G. Vinodhini, R.M. Chandrasekaran, Predicting Students' Academic Performance in the University Using Meta Decision Tree Classifiers., J. Comput. Sci. 14 (2018) 654–662.

[6] M. Bucos, B. Drăgulescu, Predicting student success using data generated in traditional educational environments, TEM Journal 7 (2018) 617.

[7] J. Xu, Y. Han, D. Marcu, M. Van Der Schaar, Progressive prediction of student performance in college programs, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2017.

[8] X. Liu, Y. Ding, H. Tang, F. Xiao, A data mining-based framework for the identification of daily electricity usage patterns and anomaly detection in building electricity consumption data, Energy Build 231 (2021) 110601.

[9] W. Villegas-Ch, S. Luján-Mora, Analysis of data mining techniques applied to LMS for personalized education, in: 2017 IEEE World Engineering Education Conference (EDUNINE), IEEE, 2017: pp. 85–89.

[10] F. Ünal, Data mining for student performance prediction in education, Data Mining-Methods, Applications and Systems 28 (2020) 423–432.

[11] L. Masangu, A. Jadhav, R. Ajoodha, Predicting student academic performance using data mining techniques, Advances in Science, Technology and Engineering Systems Journal 6 (2021) 153–163.

[12] L.D. Yulianto, A. Triayudi, I.D. Sholihati, Implementation Educational Data Mining For Analysis of Student Performance Prediction with Comparison of K-Nearest Neighbor Data Mining Method and Decision Tree C4. 5: Implementation Educational Data Mining For Analysis of Student Performance Prediction w, Jurnal Mantik 4 (2020) 441–451.

[13] S. Batool, J. Rashid, M.W. Nisar, J. Kim, H.-Y. Kwon, A. Hussain, Educational data mining to predict students' academic performance: A survey study, Educ Inf Technol (Dordr) 28 (2023) 905–971.

[14] B. Mehboob, Predicting student performance and risk analysis by using data mining approach, (2023).

[15] B. Sadaghat, G.G. Tejani, S. Kumar, Predict the Maximum Dry Density of soil based on Individual and Hybrid Methods of Machine Learning, Advances in Engineering and Intelligence Systems 2 (2023).

[16] A.P. Alfiani, F.A. Wulandari, Mapping student's performance based on data mining approach (a case study), Agriculture and Agricultural Science Procedia 3 (2015) 173–177.

[17] T. Mahboob, S. Irfan, A. Karamat, A machine learning approach for student assessment in E-learning using Quinlan's C4. 5, Naive Bayes and Random Forest algorithms, in: 2016 19th International Multi-Topic Conference (INMIC), IEEE, 2016: pp. 1–8.

[18] M. Pandey, S. Taruna, Towards the integration of multiple classifier pertaining to the Student's performance prediction, Perspect Sci (Neth) 8 (2016) 364–366.

[19] Purnawansyah, Haviluddin, K-Means clustering implementation in network traffic activities, in: 2016 International Conference on Computational Intelligence and Cybernetics, IEEE, 2016: pp. 51–54.

[20] Y.K. Saheed, T.O. Oladele, A.O. Akanni, W.M. Ibrahim, Student performance prediction based on data mining classification techniques, Nigerian Journal of Technology 37 (2018) 1087–1091.

[21] G. Deeva, J. De Smedt, C. Saint-Pierre, R. Weber, J. De Weerdt, Predicting student performance using sequence classification with time-based windows, Expert Syst Appl 209 (2022) 118182.

[22] A.D. Kumar, R.P. Selvam, V. Palanisamy, Hybrid classification algorithms for predicting student performance, in: 2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS), IEEE, 2021: pp. 1074–1079.

[23] E.H. Yossy, Y. Heryadi, Comparison of data mining classification algorithms for student performance, in: 2019 IEEE International Conference on Engineering, Technology and Education (TALE), IEEE, 2019: pp. 1–4.

[24] H. Pallathadka, A. Wenda, E. Ramirez-Asís, M. Asís-López, J. Flores-Albornoz, K. Phasinam, Classification and prediction of student performance data using various machine learning algorithms, Mater Today Proc 80 (2023) 3782–3785.

[25] C. Márquez-Vera, A. Cano, C. Romero, S. Ventura, Predicting student failure at school using genetic programming and different data mining approaches with high dimensional and imbalanced data, Applied Intelligence 38 (2013) 315–330.

[26] T. Hu, T. Song, Research on XGboost academic forecasting and analysis modelling, in: J Phys Conf Ser, IOP Publishing, 2019: p. 12091.

[27] D. Kabakchieva, Student performance prediction by using data mining classification algorithms, International Journal of Computer Science and Management Research 1 (2012) 686–690.

[28] D. Thammasiri, D. Delen, P. Meesad, N. Kasap, A critical assessment of imbalanced class distribution problem: The case of predicting freshmen student attrition, Expert Syst Appl 41 (2014) 321–330.

[29] F. Marbouti, H.A. Diefes-Dux, K. Madhavan, Models for early prediction of at-risk students in a course using standards-based grading, Comput Educ 103 (2016) 1–15.

[30] N. Dengen, E. Budiman, M. Wati, U. Hairah, Student Academic Evaluation using Naïve Bayes Classifier Algorithm, in: 2018 2nd East Indonesia Conference on Computer and Information Technology (EIConCIT), IEEE, 2018: pp. 104–107.

[31] H. Sibyan, J. Svajlenka, H. Hermawan, N. Faqih, A.N. Arrizqi, Thermal Comfort Prediction Accuracy with Machine Learning between Regression Analysis and Naïve Bayes Classifier, Sustainability 14 (2022) 15663.

[32] J.-S. Chou, D.-N. Truong, A novel metaheuristic optimizer inspired by behavior of jellyfish in ocean, Appl Math Comput 389 (2021) 125535.

[33] A. Alam, P. Verma, M. Tariq, A. Sarwar, B. Alamri, N. Zahra, S. Urooj, Jellyfish search optimization algorithm for mpp tracking of pv system, Sustainability 13 (2021) 11736.

[34] L. Wang, Q. Cao, Z. Zhang, S. Mirjalili, W. Zhao, Artificial rabbits optimization: A new bio-inspired meta-heuristic algorithm for solving engineering optimization problems, Eng Appl Artif Intell 114 (2022) 105082.

[35] G.S. Linoff, M.J.A. Berry, Data mining techniques: for marketing, sales, and customer relationship management, John Wiley & Sons, 2011.

[36] T. Hastie, R. Tibshirani, J.H. Friedman, J.H. Friedman, The elements of statistical learning: data mining, inference, and prediction, Springer, 2009.

[37] I.H. Witten, E. Frank, Data mining: practical machine learning tools and techniques with Java implementations, Acm Sigmod Record 31 (2002) 76–77.

[38] X. Luo, Efficient English text classification using selected machine learning techniques, Alexandria Engineering Journal 60 (2021) 3401–3409.

[39] Q.A. Al-Radaideh, A. Al Ananbeh, E. Al-Shawakfa, A classification model for predicting the suitable study track for school students, Int. J. Res. Rev. Appl. Sci 8 (2011) 247–252.

[40] R.R. Kabra, R.S. Bichkar, Performance prediction of engineering students using decision trees, Int J Comput Appl 36 (2011) 8–12.

[41] C. Márquez-Vera, A. Cano, C. Romero, S. Ventura, Predicting student failure at school using genetic programming and different data mining approaches with high dimensional and imbalanced data, Applied Intelligence 38 (2013) 315–330.

[42] D. Kabakchieva, Student performance prediction by using data mining classification algorithms, International Journal of Computer Science and Management Research 1 (2012) 686–690.

[43] N.T. Nghe, P. Janecek, P. Haddawy, A comparative analysis of techniques for predicting academic performance, in: 2007 37th Annual Frontiers in Education Conference-Global Engineering: Knowledge without Borders, Opportunities without Passports, IEEE, 2007: pp. T2G-7.

[44] E. Osmanbegovic, M. Suljic, Data mining approach for predicting student performance, Economic Review: Journal of Economics and Business 10 (2012) 3–12.

# A Deep Learning-based Framework for Vehicle License Plate Detection

Deming Yang[1*], Ling Yang[2]

Huanghai College of Automotive Engineering, Liaoning Jidian Polytechnic, Dandong 118009, Liaoning, China
Department of Information Engineering, Liaoning Jidian Polytechnic, Dandong 118009, Liaoning, China

*Abstract*—In the contemporary landscape of smart transportation systems, the imperative role of intelligent traffic monitoring in bolstering efficiency, safety, and sustainability cannot be overstated. Leveraging recent strides in computer vision, machine learning, and data analytics, this study addresses the pressing need for advancements in car license plate recognition within these systems. Employing an innovative approach based on the YOLOv5 architecture in deep learning, the study focuses on refining the accuracy of license plate recognition. A bespoke dataset is meticulously curated to facilitate a comprehensive evaluation of the proposed methodology, with extensive experiments conducted and metrics such as precision, recall, and F1-score employed for assessment. The outcomes underscore the efficacy of the approach in significantly enhancing the precision and accuracy of license plate recognition using performance evaluation of the proposed method. This tailored dataset ensures a rigorous evaluation, affirming the practical viability of the proposed approach in real-world scenarios. The study not only showcases the successful application of deep learning and YOLOv5 in achieving accurate license plate detection and recognition but also contributes to the broader discourse on advancing intelligent traffic monitoring for more robust and efficient smart transportation systems.

*Keywords*—*Intelligent traffic monitoring; smart transportation; deep learning; Yolov5; performance evaluation*

## I. INTRODUCTION

Intelligent Traffic Monitoring plays a vital role in the development of smart transportation systems, aiming to enhance the efficiency, safety, and sustainability of modern urban mobility [1], [2]. By utilizing advanced technologies and data-driven approaches, intelligent traffic monitoring enables real-time analysis, prediction, and management of traffic conditions [3], [4]. It encompasses various components such as sensor networks, data integration, traffic analysis algorithms, and intelligent decision-making systems [5]. Efficient traffic monitoring is essential for optimizing traffic flow, reducing congestion, and improving overall transportation experiences.

Recent years have witnessed significant advancements in intelligent traffic monitoring technologies, driven by the rapid progress in computer vision, machine learning, and data analytics[6], [7]. Fig. 1 demonstrates the schematic of an intelligent transportation system [8]. These advancements have enabled the development of sophisticated systems that can automatically capture, process, and analyze vast amounts of traffic data in real-time. Technologies like video surveillance cameras, radar systems, as well as connected vehicles contribute to the data collection process [9]. Meanwhile, advanced algorithms and analytics techniques, including deep learning [10], have emerged as powerful tools for traffic analysis, pattern recognition, and prediction.

The research significance of intelligent traffic monitoring in smart transportation is paramount. Accurate and efficient monitoring systems are crucial for traffic management authorities, urban planners, and policymakers to make informed decisions regarding infrastructure planning, traffic control strategies, and resource allocation. Furthermore, intelligent traffic monitoring has the potential to improve safety by enabling early detection of traffic incidents and facilitating timely emergency responses. It also contributes to reducing energy consumption, mitigating environmental impacts, and enhancing overall transportation system resilience.

Among the various technologies used in intelligent traffic monitoring [11], vision-based systems have attracted considerable attention from researchers [12]. Vision-based systems use computer vision techniques to extract relevant information from visual data, like images or videos captured by surveillance cameras. The appeal of vision-based approaches lies in their ability to provide rich and detailed information about traffic conditions, including vehicle movements, license plate recognition [13], and traffic flow analysis. These systems have been continuously evolving, with the introduction of more sophisticated computer vision algorithms and deep learning-based methods.

Several studies have focused on computer vision-based intelligent traffic monitoring systems. These studies have explored different aspects, such as object detection, tracking, license plate recognition, and traffic flow estimation [14], [15]. Fig. 2 shows the schematic of a license plate recognition system. Recently, deep learning-based approaches [16], particularly those employing You Only Look Once (YOLO) architecture [17], have gained attention due to their superior performance in object detection tasks [18], [19,30]. These methods leverage large-scale datasets and powerful deep neural network architectures to achieve high accuracy as well as real-time processing capabilities [20]. However, despite these advancements, there still exist certain limitations and research gaps that need to be addressed to enhance further the effectiveness and robustness of intelligent traffic monitoring systems.

Fig. 1. Schematic of an intelligent transportation system.



Fig. 2. Schematic of a license plate recognition system [21].

This study is motivated by the essential role of Intelligent Traffic Monitoring in advancing smart transportation systems for enhanced efficiency, safety, and sustainability in urban mobility. Recent technological advancements, particularly in computer vision and machine learning, have enabled sophisticated monitoring systems, but existing approaches still have limitations. The focus is on vision-based systems, utilizing computer vision techniques for detailed traffic information, including license plate recognition. To address research gaps, the study proposes a novel approach using deep learning with the YOLOv5 architecture to improve car license plate recognition. The aim is to contribute to the field by conducting extensive experiments with a custom dataset, encompassing training, validation, and testing processes, to enhance the effectiveness and robustness of intelligent traffic monitoring systems.

This study presents a method that uses a deep learning framework with the YOLOv5 architecture to improve car license plate recognition. To evaluate the effectiveness of the proposed approach, generating a custom dataset and conduct extensive experiments involving training, validation, and testing processes.

The following is a list of the major research contributions made in this article:

*1)* Improving the accuracy of car license plate recognition by proposing a novel approach based on a deep learning technique with the YOLOv5 architecture, resulting in superior recognition performance compared to existing methods.

*2)* Developing and utilizing a custom dataset encompassing diverse car license plate images to facilitate comprehensive evaluation and benchmarking of the proposed YOLOv5-based approach for license plate recognition.

*3)* Conducting extensive experiments to assess the effectiveness as well as the performance of the suggested YOLOv5-based approach, including training, validation, and testing processes, and evaluating metrics regarding F1-score, recall and precision performance metrics.

The remainder of this article includes the following parts: Section I presents the introduction. The literature review is discussed in Section II. Section III presents the material and method. Experimental results are discussed in Section IV. Finally, Section V concludes the paper.

## II. Related Works

Wang et al. [22] proposed a light convolutional neural network (CNN) approach for end-to-end car license plate recognition and detection. The method involves training a lightweight CNN model that can directly detect and recognize license plates in images. The key features of the approach include its simplicity, efficiency, and ability to handle various license plate designs and environmental conditions. The findings show that the proposed light CNN achieves competitive performance regarding accuracy and processing speed, outperforming traditional methods. However, the limitation of the method lies in its dependence on well-labeled training data and potential challenges in handling extremely blurred or distorted license plate images. Further research is required to address these limitations and enhance the system's adaptability to complex real-world scenarios.

Pustokhina et al. [23] developed an automatic vehicle license plate recognition system for intelligent transportation systems, utilizing optimal K-means clustering in combination with a convolutional neural network (CNN). The method demonstrates high accuracy and efficiency in recognizing license plate characters, leveraging K-means clustering for region extraction and a CNN for character segmentation and recognition. The system achieves competitive performance compared to existing approaches and proves effective in diverse scenarios. However, the paper acknowledges limitations related to region extraction accuracy and adaptability to varying license plate designs and environmental conditions, suggesting further research for enhancing system robustness.

He et al. [24] presented a robust automatic recognition system for Chinese license plates in natural scenes. The proposed method employs a combination of image preprocessing, character segmentation, and recognition algorithms. The key features include adaptive thresholding, connected component analysis, as well as a CNN-based classifier for character recognition. The system achieves high accuracy in license plate recognition, even in challenging scenarios with variations in lighting conditions, plate orientations, and background clutter. The findings show that the suggested system outperforms existing methods on benchmark datasets, demonstrating its robustness and effectiveness. However, the limitations of the system lie in its sensitivity to complex backgrounds and the requirement for well-segmented license plate regions.

Kaur et al. [25] presented an automatic license plate recognition system for vehicles using a convolutional neural network (CNN). The method involves training the CNN on a large dataset of license plate images to learn the patterns and features necessary for accurate recognition. Key features of the system include the ability to extract license plate regions, segment individual characters, and recognize them using the trained CNN. The findings illustrate that the suggested CNN-based approach achieves high accuracy in license plate recognition, outperforming traditional methods. However, the system's limitation lies in its dependence on well-labeled training data and its susceptibility to variations in license plate designs, environmental conditions, and image quality.

A low-cost Internet of Things (IoT) based Arabic license plate recognition model for smart parking systems presented by Abdellatif et al. [26]. The method uses a combination of machine learning algorithms and image processing techniques to detect as well as recognize Arabic license plates. Key features of the approach include its cost-effectiveness, reliance on IoT infrastructure, as well as the ability to recognize Arabic license plates accurately. The findings demonstrate the effectiveness of the suggested model in real-world scenarios, achieving high accuracy in license plate recognition for smart parking systems. However, the limitation of the system lies in its focus on Arabic license plates, limiting its applicability to other regions with different license plate formats. Further research is recommended to explore the model's generalizability to different languages and license plate designs.

Shi et al. [27] proposed a License Plate Recognition System (LPRS) that combines an improved version of the YOLOv5 object detection algorithm and a GRU model for accurate license plate recognition and detection. The method involves two stages: license plate detection using the improved YOLOv5, which incorporates feature fusion and attention mechanism to enhance the detection performance, and license plate recognition using a GRU model, which learns the sequential patterns of characters on the license plate. The key features of the proposed system include improved accuracy in license plate recognition and detection, robustness to varying lighting and environmental conditions, and efficient processing speed. The experimental outcomes demonstrate superior performance compared to existing methods, achieving high accuracy rates in license plate recognition and detection tasks. However, the limitation of the system lies in its reliance on a large dataset for training, which may pose challenges in scenarios with limited data availability.

Despite the notable advancements in vehicle license plate recognition (VLPR) systems, a research gap persists in

achieving a more accurate and robust method for license plate detection. Although existing studies, have made significant contributions by introducing various approaches combining the CNNs, clustering techniques, and innovative algorithms for license plate detection and recognition. However, common limitations across these studies include dependencies on well-labeled training data, challenges in handling blurred or distorted license plate images, adaptability issues to varying designs and environmental conditions, and sensitivity to complex backgrounds. While these studies demonstrate promising results, the need for a more accurate and adaptable VLPR system that addresses these limitations is evident. Therefore, the research gap lies in developing a method that not only surpasses current accuracy rates but also ensures robustness in challenging real-world scenarios, such as scenarios with limited training data availability or diverse license plate designs and environmental conditions. Closing this gap is crucial for advancing the field and facilitating the deployment of more reliable and effective vehicle license plate detection systems in smart transportation and intelligent traffic monitoring applications.

## III. METHODOLOGY

In this study, the researchers propose an algorithm for license plate recognition that is based on the YOLOv5 model. The algorithm leverages the strengths of YOLOv5, which is known for its real-time object detection capabilities, to recognize license plates in images or videos effectively. By using YOLOv5 inspired by [17], [28], [29], the algorithm aims to achieve accurate as well as efficient license plate recognition, offering potential applications in various domains, such as traffic monitoring for smart transportation applications.

In this study, a method for license plate recognition is proposed, capitalizing on the strengths of the YOLOv5 model renowned for its real-time object detection capabilities. The algorithm aims to achieve both accuracy and efficiency in license plate recognition, with potential applications in diverse domains, particularly in smart transportation for traffic monitoring. The proposed methodology involves a multi-step process for model generation within the YOLOv5 framework. The first step, Data Preparation, entails collecting a labeled dataset of license plate images, annotating license plate regions, and assigning corresponding labels. Model Training follows, where the YOLOv5 model is trained on the annotated dataset, optimizing parameters for accurate license plate detection. The subsequent step, Model Evaluation, utilizes metrics like mean average precision (mAP) and detection accuracy to assess the trained model's performance. Inference involves utilizing the trained YOLOv5 model for license plate detection on unseen images or videos. The final step, post-processing, employs techniques such as non-maximum suppression (NMS) to refine predictions and ensure the most confident and accurate license plate detections are retained. Noteworthy efforts are invested in Data Preparation, which includes image annotation and data augmentation for creating a custom dataset that enhances the model's robustness and generalization capabilities. Model Training involves splitting the dataset, training the YOLO model, and fine-tuning based on the validation set. The proposed algorithm provides a comprehensive approach to license plate recognition,

addressing real-world challenges and demonstrating potential advancements in the field.

The following categories apply to the steps involved in model generation in YOLOv5:

*1) Data preparation:* Gather a labeled dataset of license plate images. Annotate the license plate regions with bounding boxes and assign corresponding labels.

*2) Model training:* Train the YOLOv5 model on the annotated license plate dataset. This involves optimizing the model's parameters to detect license plates accurately.

*3) Model evaluation:* Utilize evaluation metrics like mean average precision (mAP) and detection accuracy to assess the performance of the trained model. This step helps determine the model's effectiveness in license plate recognition.

*4) Inference:* Utilize the trained YOLOv5 model to perform license plate detection on unseen images or videos. The model identifies the license plate regions and provides the corresponding bounding boxes and class predictions.

*5) Post-processing:* Apply post-processing methods like non-maximum suppression (NMS) to filter out redundant bounding boxes and retain the most accurate license plate detections.

### B. Data Preparation

Preparing data for training a YOLOv5 model using collected images involves two key steps: image annotation and data augmentation. This study collected images from internet resources. In image annotation, the specific objects of interest, such as license plates, are manually labeled by drawing bounding boxes around them and assigning corresponding class labels. This step ensures that the model learns to detect the desired objects accurately. Data augmentation, on the other hand, involves applying various transformations to the images to expand the training dataset artificially. Techniques like cropping, rotation, scaling, flipping, color jittering, and noise injection are used to simulate real-world variations, improving the model's ability to generalize to various scenarios.

By combining image annotation and data augmentation, a custom dataset is created for training a YOLOv5 model. The annotated images provide the necessary ground truth information, enabling the model to learn the spatial location and class labels of license plates. Data augmentation diversifies the dataset, introducing variations in object appearance, scale, and environmental conditions. This helps the model become more robust and generalize well to unseen license plate images. These steps prepare the data as well as ensure the effectiveness and accuracy of the YOLOv5 model for license plate detection tasks.

### C. Model Training

To generate a YOLOv5 model using a custom dataset and split it into training, testing, and validation sets, you can follow these steps:

*1) Dataset split:* First, split the custom dataset into three subsets: testing, validation, also training. Allocate 75% of the data for training, 15% for validation, also the remaining 10%

for testing. Ensure that the splits are representative and maintain a balanced distribution of license plate images across all subsets.

*2) Training module:* The training module involves training the YOLO model on the training dataset. During training, the model learns to detect license plates by optimizing its parameters utilizing a training algorithm like Adam optimizer or stochastic gradient descent (SGD). The training process iteratively updates the model's weights based on the loss function, which consists of localization loss (IoU loss) as well as confidence loss (objectness loss). The model is exposed to the training dataset, and backpropagation is used to update the weights, gradually improving the model's accuracy in detecting license plates.

*3) Validation module:* The validation module is used to assess the model's performance and tune its hyperparameters during training. The validation set, comprising 15% of the custom dataset, is utilized to evaluate the model's accuracy and generalization ability. The trained model is run on the validation set, and metrics like mean average precision (mAP) are calculated to measure the quality of license plate detection. To optimize the model's performance, adjustments to the training procedure, model architecture, or hyperparameters might be made in light of the validation results.

*4) Testing module:* The testing module is the final stage, where the trained YOLO model is appraised on the testing dataset, which also consists of 10% of the custom dataset. The testing dataset contains unseen license plate images that were not used during training or validation. The model's performance is assessed by running it on the testing dataset, also measuring metrics like mAP, detection accuracy, and false positives/negatives. This module provides an understanding of the model's real-world performance as well as its ability to recognize license plates in unseen and new scenarios accurately.

*D. Model Evaluation*

Assessing the model is vital to validate the proficiency of a YOLOv5 model for recognizing license plates. It provides a means to objectively assess the model's performance, measure its accuracy and generalization ability, and identify areas for improvement. By evaluating the model on independent datasets, stakeholders can gain confidence in its effectiveness, ensure it meets desired requirements, and make informed decisions regarding its deployment and usage. To evaluate a generated YOLOv5 model for license plate recognition using performance metrics like F1-score, recall, and precision, the following steps are:

*1) Intersection over Union (IoU) calculation:* For each predicted bounding box, calculate the Intersection over Union (IoU) with the corresponding ground truth bounding box. IoU examines the overlap between the predicted and ground truth bounding boxes to determine if detection is a true positive or a false positive.

*2) Precision:* Precision measures the proportion of correctly detected license plates among all the predicted license plates. It is computed as the ratio of true positives (correct detections) to the sum of false positives as well as true positives (incorrect detections).

*3) Recall:* Recall measures the proportion of correctly detected license plates among all the ground truth license plates. It is computed as the ratio of true positives to the total of false negatives also true positives (missed detections).

*4) F1-score:* The F1-score, which provides a balanced assessment of the model's performance, is the harmonic mean of precision and recall. It is computed as 2 * ((precision * recall) / (precision + recall)).

Examine the recall, precision, and F1-score values to assess the model's performance in license plate recognition. Greater precision signifies fewer false positive detections, whereas higher recall signifies fewer false negatives. The F1-score provides an overall evaluation by considering both precision and recall. The details of performance metrics results are presented in the following section.

*E. Inference*

Inference refers to the process of utilizing a trained YOLOv5 model to perform license plate detection on new, unseen images or videos. Once the YOLOv5 model has been trained on a dataset, it learns to recognize license plates by detecting and localizing them within an image. During inference, the trained model takes an input image or frame from a video and processes it through the model's architecture.

The inference process involves passing the input image through the YOLOv5 model, which applies a series of convolutional layers, down-sampling, and feature extraction to identify potential license plate regions. The model predicts bounding box coordinates as well as class probabilities for each detected object, including license plates. These predictions are generated based on the learned patterns and features acquired during the training phase. After running inference, the YOLOv5 model provides the predicted bounding box coordinates and class labels for detected license plates. This information allows for precise localization of license plates within the image or video frame. The model's output can be further processed to extract the license plate region and perform subsequent tasks such as character recognition or vehicle identification.

*F. Post-processing*

In the post-processing stage of license plate recognition using a trained YOLOv5 model, techniques such as non-maximum suppression (NMS) are applied to refine the model's predictions. NMS helps eliminate redundant bounding boxes by retaining only the most confident and accurate detections. By comparing the confidence scores and utilizing intersection over union (IoU) calculations, NMS suppresses overlapping detections and ensures that only the highest-scoring detection for each object is retained. This post-processing step improves the quality and accuracy of license plate recognition results by removing duplicate detections and producing cleaner outputs.

## IV. RESULTS AND PERFORMANCE EVALUATION

Performance evaluation and experimental outcomes for a YOLOv5 model in license plate recognition involve assessing the model's performance utilizing metrics like accuracy, recall, precision, F1-score and mAP. Fig. 3 shows a sample result of the license plate recognition. The model is tested on a separate dataset with ground truth annotations, and its predictions are compared against these annotations to calculate the metrics. The results provide quantitative measures of the model's ability to detect and localize license plates accurately.

The analysis of experimental results allows for an evaluation of the YOLOv5 model's overall performance and helps identify areas for improvement. By examining metrics like precision, accuracy, recall, and F1-score, the model's effectiveness in license plate recognition can be assessed. The mAP metric provides insights into the trade-off between precision and recall, while IoU analysis helps gauge the localization accuracy. Comparisons with other models or benchmarks further contribute to understanding the model's relative strengths and weaknesses. These findings guide iterative improvement, enabling researchers to refine the YOLOv5 model's architecture, training process, or other parameters to enhance its license plate recognition capabilities.

For the license plate recognition using the YOLOv5 model, performance metrics such as precision are commonly employed to assess the model's accuracy. For this evaluation, following performance metrics are used,

True Positive (TP): It refers to the cases where the YOLOV5 model correctly predicts a license plate as positive (i.e., correctly identifying a license plate when one exists).

True Negative (TN): It represents the cases where the YOLOV5 model correctly predicts the absence of a license plate as negative (i.e., correctly identifying the absence of a license plate when none exists).

False Positive (FP): It occurs when the YOLOV5 model incorrectly predicts a license plate when there isn't one (i.e., incorrectly identifying a license plate when none exists).

False Negative (FN): It refers to the cases where the YOLOV5 model fails to predict a license plate when one exists (i.e., incorrectly not identifying a license plate when one exists).

Based on these definitions, the ability of a model to accurately identify positive predictions is measured by its precision. It is calculated using the following formula:

$$Precision = TP/(TP + FP)$$

Furthermore, the recall measures the proportion of correctly predicted license plates out of all actual license plates. It is calculated using the following formula:

$$Recall = TP/(TP + FN)$$

Precision quantifies the proportion of correctly predicted license plates out of all predicted license plates. A higher precision demonstrates fewer false positives, meaning that the model is more accurate in identifying true license plates. The recall quantifies the model's ability to recognize all positive instances correctly. A higher recall demonstrates fewer false negatives, indicating that the model can successfully detect most of the license plates present. Fig. 4 and Fig. 5 show the result of precision and recall performance metrics.

As displayed in Fig. 4, precision is a performance metric utilized to appraise the accuracy of a YOLOv5 model in license plate recognition. It measures the proportion of correctly detected license plates among the predicted ones. In precision evaluation, the Y-axis demonstrates precision values, while the X-axis represents confidence values or thresholds. The precision curve visualizes the trade-off between precision as well as recall at various confidence thresholds, helping determine the optimal threshold that balances precision and recall. Higher precision indicates fewer false positives, making it important for tasks like license plate recognition in real-world scenarios.

As shown in Fig. 5, in the recall evaluation, the Y-axis denotes the recall values, while the X-axis denotes the confidence values or the threshold applied to the model's predictions. To calculate recall, the forecasted bounding boxes, as well as their associated confidence scores, are sorted in descending order. A threshold is set on the confidence scores to determine which detections are regarded as positive forecasts. By varying the threshold, different levels of confidence are required for a detection to be considered positive. The recall is then computed as the ratio of true positives (correctly detected license plates) to the total of true positives also false negatives (missed detections).



Fig. 3. Samples result of the license plate recognition.

Fig. 4. The result of precision metric.



Fig. 5. Result of recall metric.

Moreover, as shown in Fig. 5, the recall curve is a graphical representation of recall as a function of the confidence threshold. The curve is generated by plotting the recall values at different confidence thresholds on the Y-axis against the corresponding confidence values on the X-axis. As the threshold decreases (lower confidence threshold), more detections are considered positive, resulting in increased recall but potentially lower precision. The recall curve gives insights into the model's ability to detect license plates across different confidence thresholds.

The F1-score is the harmonic mean of a recall and precision and is computed utilizing the following formula:

$$F1 - score = 2 * (precision * recall)/(precision + recall)$$

The F1-score curve and how it relates to the YOLOv5 model evaluation for license plate recognition. As displayed in Fig. 6, the Y-axis of the F1-score curve represents the F1-score itself. It is a continuous range of values between 0 and 1, where 1 represents a perfect F1-score (indicating high precision and recall), and 0 represents the worst score (indicating poor performance in both precision and recall). The X-axis of the F1-score curve typically demonstrates the confidence values or the decision thresholds used by the YOLOv5 model. During inference, the YOLOv5 model predicts bounding boxes for license plates along with their confidence scores, indicating how confident the model is in its predictions. The confidence threshold is a value used to determine whether a predicted bounding box and its associated label (license plate) are considered valid or not. The F1-score curve is obtained by plotting the F1-scores at different confidence thresholds. By varying the threshold, the model's precision and recall trade-off can be analyzed. Generally, as the confidence threshold increases, the model becomes more conservative in its predictions, resulting in higher precision but potentially lower recall. Conversely, lowering the threshold may lead to higher recall but lower precision. The F1-score curve provides insights into the model's performance at various confidence thresholds, allowing you to choose a suitable threshold based on ythe specific requirements.



Fig. 6. Result of F1-score metric



Fig. 7. Result of precision-recall (PR) curve

As shown in Fig. 7, the precision-recall curve is a graphical representation utilized to measure the performance of the YOLOv5 model in license plate recognition tasks. The Y-axis represents precision values ranging from 0 to 1, where higher values display better accuracy in positive predictions. The X-axis represents recall values, also ranging from 0 to 1, where higher values indicate a higher ability to detect positive instances. The precision-recall curve for the YOLOv5 model evaluation in license plate recognition showcases how precision and recall vary with changing confidence thresholds. Adjusting the threshold allows for different trade-offs between recall and precision. Rising the threshold tends to increase precision but may result in lower recall, while decreasing the threshold may raise recall at the expense of precision. By analyzing the curve, users can evaluate the model's performance at different thresholds as well as select the appropriate threshold based on their priorities. Whether prioritizing precision to minimize false positives or recall to minimize false negatives, the precision-recall curve offers a visual representation to assess and fine-tune the YOLOv5 model's performance for license plate recognition tasks.

The discussion of the presented results involves a comprehensive evaluation of the YOLOv5 model's performance in license plate recognition, utilizing various metrics such as accuracy, recall, precision, F1-score, and mAP. The model undergoes testing on a separate dataset with ground truth annotations, and the quantitative measures obtained offer insights into its ability to accurately detect and localize license plates. The precision metric is utilized to assess the model's accuracy, considering true positives (correctly predicted license plates), true negatives (correctly predicted absence of license plates), false positives (incorrect predictions of license plates), and false negatives (missed predictions). The precision curve, as illustrated in Fig. 4, demonstrates the trade-off between precision and recall at different confidence thresholds, aiding in determining an optimal threshold for balancing accuracy. Similarly, the recall metric, depicted in Figure 5, measures the model's ability to correctly predict license plates and is analyzed across various confidence thresholds. The recall curve offers insights into the model's performance at different thresholds, showcasing the trade-off between recall and precision. The F1-score, presented in Fig. 6, represents the harmonic mean of precision and recall, providing a balanced assessment of the model's performance across different confidence thresholds. Finally, the precision-recall (PR) curve, as shown in Fig. 7, visually represents the model's performance in license plate recognition, enabling the analysis of trade-offs between precision and recall at different confidence thresholds. These results not only contribute to a comprehensive understanding of the YOLOv5 model's strengths and weaknesses but also guide potential refinements in its architecture, training processes, and parameters to enhance license plate recognition capabilities.

## V. Conclusion

This study introduces a novel and improved strategy for car license plate recognition, leveraging the YOLOv5 architecture within a deep learning framework. The proposed approach demonstrates a commitment to achieving superior performance compared to existing methods, with a particular focus on enhancing accuracy and efficiency. The development of a custom dataset and the thorough evaluation through extensive experiments, utilizing metrics like recall, precision, and F1-score, substantiate the efficacy of the proposed methodology. Despite the promising outcomes, it is essential to acknowledge certain limitations. The dependence on well-labeled training data and potential challenges in handling extremely blurred or distorted license plate images pose constraints on the adaptability of the approach to complex real-world scenarios. Future research endeavors could address these limitations by exploring innovative techniques for handling diverse and challenging environmental conditions. Additionally, an exciting avenue for further exploration lies in integrating multiple sensors for comprehensive data collection, potentially leading to the development of real-time decision-making systems for traffic control and optimization. This study, while providing a valuable contribution to the field, sets the stage for future investigations aimed at advancing intelligent traffic monitoring systems and ultimately contributing to the ongoing evolution of smart transportation. Furthermore, this study acknowledges the limitations of the proposed approach, particularly its reliance on well-labeled training data and potential challenges in handling distorted license plate images, future research could delve into the development of more robust algorithms capable of addressing these complexities. Exploring advanced image processing techniques or incorporating additional pre-processing steps may enhance the adaptability of the system to varying real-world scenarios. Moreover, there is a promising avenue for research in the integration of multiple sensors, such as cameras, LIDAR, and radar, to create a more comprehensive and diverse dataset for training. This holistic approach could significantly contribute to the system's ability to handle different environmental conditions and improve generalization across various scenarios. Additionally, the realization of real-time decision-making systems for traffic control and optimization holds great potential, necessitating research efforts to design and implement algorithms [30] capable of providing instantaneous responses to dynamic traffic conditions. By addressing these research directions, future studies can contribute to the ongoing evolution of intelligent traffic monitoring systems, ensuring their robustness, adaptability, and efficiency in the ever-changing landscape of smart transportation.

## References

[1] R. Li, S. Wang, P. Jiao, and S. Lin, "Traffic control optimization strategy based on license plate recognition data," Journal of traffic and transportation engineering (English edition), vol. 10, no. 1, pp. 45–57, 2023.

[2] N.-A.- Alam, M. Ahsan, M. A. Based, and J. Haider, "Intelligent system for vehicles number plate detection and recognition using convolutional neural networks," Technologies (Basel), vol. 9, no. 1, p. 9, 2021.

[3] V. S. R. Kosuru, A. K. Venkitaraman, V. D. Chaudhari, N. Garg, A. Rao, and A. Deepak, "Automatic Identification of Vehicles in Traffic using Smart Cameras," in 2022 5th International Conference on Contemporary Computing and Informatics (IC3I), IEEE, 2022, pp. 1009–1014.

[4] N. A. Khan, N. Z. Jhanjhi, S. N. Brohi, R. S. A. Usmani, and A. Nayyar, "Smart traffic monitoring system using unmanned aerial vehicles (UAVs)," Comput Commun, vol. 157, pp. 434–443, 2020.

[5] N. do V. Dalarmelina, M. A. Teixeira, and R. I. Meneguette, "A real-time automatic plate recognition system based on optical character recognition and wireless sensor networks for ITS," Sensors, vol. 20, no. 1, p. 55, 2019.

[6] S. Nayak and K. Katakiya, "Machine Vision Based Intelligent Traffic Management Tool," IJRAR-International Journal of Research and Analytical Reviews (IJRAR), vol. 6, no. 2, pp. 83–88, 2019.

[7] J. Tang, L. Wan, J. Schooling, P. Zhao, J. Chen, and S. Wei, "Automatic number plate recognition (ANPR) in smart cities: A systematic review on technological advancements and application cases," Cities, vol. 129, p. 103833, 2022.

[8] S. Cao, Y.-J. Wu, and F. Jin, "New Radar Sensor Technology for Intelligent Multimodal Traffic Monitoring at Intersections," 2021.

[9] K. T. Islam et al., "A vision-based machine learning method for barrier access control using vehicle license plate authentication," Sensors, vol. 20, no. 12, p. 3578, 2020.

[10] R. Balia, S. Barra, S. Carta, G. Fenu, A. S. Podda, and N. Sansoni, "A deep learning solution for integrated traffic control through automatic license plate recognition," in Computational Science and Its Applications–ICCSA 2021: 21st International Conference, Cagliari, Italy, September 13–16, 2021, Proceedings, Part III 21, Springer, 2021, pp. 211–226.

[11] H. O. Al-Sakran, "Intelligent traffic information system based on integration of Internet of Things and Agent technology," Int J Adv Comput Sci Appl, vol. 6, no. 2, pp. 37–43, 2015.

[12] Lubna, N. Mufti, and S. A. A. Shah, "Automatic number plate Recognition: A detailed survey of relevant algorithms," Sensors, vol. 21, no. 9, p. 3028, 2021.

[13] S.-O. Shim, R. Imtiaz, A. Siddiq, and I. R. Khan, "License Plates Detection and Recognition with Multi-Exposure Images," International Journal of Advanced Computer Science and Applications, vol. 13, no. 4, 2022.

[14] F. Zantalis, G. Koulouras, S. Karabetsos, and D. Kandris, "A review of machine learning and IoT in smart transportation," Future Internet, vol. 11, no. 4, p. 94, 2019.

[15] H. Song, H. Liang, H. Li, Z. Dai, and X. Yun, "Vision-based vehicle detection and counting system using deep learning in highway scenes," European Transport Research Review, vol. 11, no. 1, pp. 1–16, 2019.

[16] B. S. Abunasser, M. R. J. AL-Hiealy, A. M. Barhoom, A. R. Almasri, and S. S. Abu-Naser, "Prediction of instructor performance using machine and deep learning techniques," International Journal of Advanced Computer Science and Applications, vol. 13, no. 7, 2022.

[17] A. R. Youssef, A. A. Ali, and F. R. Sayed, "Real-time Egyptian License Plate Detection and Recognition using YOLO," International Journal of Advanced Computer Science and Applications, vol. 13, no. 7, 2022.

[18] D. J. I. Z. Chen, "Automatic vehicle license plate detection using K-means clustering algorithm and CNN," Journal of Electrical Engineering and Automation, vol. 3, no. 1, pp. 15–23, 2021.

[19] R.-C. Chen, "Automatic License Plate Recognition via sliding-window darknet-YOLO deep learning," Image Vis Comput, vol. 87, pp. 47–56, 2019.

[20] S. N. H. S. Abdullah, K. Omar, A. S. Zaini, M. Petrou, and M. Khalid, "Determining adaptive thresholds for image segmentation for a license plate recognition system," International Journal of Advanced Computer Science and Applications, vol. 7, no. 6, 2016.

[21] D. Harjani, M. Jethwani, N. Keswaney, and S. Jacob, "Automated parking management system using license plate recognition," International Journal of Computer Technology and Applications, vol. 4, no. 5, p. 741, 2013.

[22] W. Wang, J. Yang, M. Chen, and P. Wang, "A light CNN for end-to-end car license plates detection and recognition," IEEE Access, vol. 7, pp. 173875–173883, 2019.

[23] I. V. Pustokhina et al., "Automatic vehicle license plate recognition using optimal K-means with convolutional neural network for intelligent transportation systems," Ieee Access, vol. 8, pp. 92907–92917, 2020.

[24] M.-X. He and P. Hao, "Robust automatic recognition of Chinese license plates in natural scenes," Ieee Access, vol. 8, pp. 173804–173814, 2020.

[25] P. Kaur, Y. Kumar, S. Ahmed, A. Alhumam, R. Singla, and M. F. Ijaz, "Automatic License Plate Recognition System for Vehicles Using a CNN.," Computers, Materials & Continua, vol. 71, no. 1, 2022.

[26] M. M. Abdellatif, N. H. Elshabasy, A. E. Elashmawy, and M. AbdelRaheem, "A low cost IoT-based Arabic license plate recognition model for smart parking systems," Ain Shams Engineering Journal, vol. 14, no. 6, p. 102178, 2023.

[27] H. Shi and D. Zhao, "License Plate Recognition System Based on Improved YOLOv5 and GRU," IEEE Access, vol. 11, pp. 10429–10439, 2023.

[28] G. Jocher et al., "ultralytics/yolov5: v5. 0-YOLOv5-P6 1280 models, AWS, Supervise. ly and YouTube integrations," Zenodo, 2021.

[29] S. Raj, Y. Gupta, and R. Malhotra, "License plate recognition system using yolov5 and cnn," in 2022 8th International Conference on Advanced Computing and Communication Systems (ICACCS), IEEE, 2022, pp. 372–377.

[30] Bangyal, W.H., Nisar, K., Soomro, T.R., Ag Ibrahim, A.A., Mallah, G.A., Hassan, N.U. and Rehman, N.U., 2022. "An Improved Particle Swarm Optimization Algorithm for Data Classification". Applied Sciences, 13(1), p.283.

# Estimation of Heating Load Consumption in Residual Buildings using Optimized Regression Models Based on Support Vector Machine

Chao WANG[1], Xuehui QIU[2]

Qinhuangdao Vocational and Technical College, Qinhuangdao, 066100, China[1]
School of Urban Geology and Engineering, Hebei GEO University, Shijiazhuang 050000, China[2]

*Abstract*—Accurate energy consumption forecasting and assessing retrofit options are vital for energy conservation and emissions reduction. Predicting building energy usage is complex due to factors like building attributes, energy systems, weather conditions, and occupant behavior. Extensive research has led to diverse methods and tools for estimating building energy performance, including physics-based simulations. However, accurate simulations often require detailed data and vary based on modeling sophistication. The growing availability of public building energy data offers opportunities for applying machine learning to predict building energy performance. This study evaluates Support Vector Regression ($SVR$) models for estimating building heating load consumption. These models encompass a single model, one optimized with the Transit Search Optimization Algorithm (TSO) and another optimized with the Coot optimization algorithm (COA). The training dataset consists of 70% of the data, which incorporates eight input variables related to the geometric and glazing characteristics of the buildings. Following the validation of 15% of the dataset, the performance of the remaining 15% is evaluated using five different assessment metrics. Among the three candidate models, Support Vector Regression optimized with the Coot optimization algorithm (SVCO) demonstrates remarkable accuracy and stability, reducing prediction errors by an average of 20% to over 50% compared to the other two models and achieving a maximum $R^2$ value of 0.992 for heating load prediction.

*Keywords—Heating load demand; prediction models; building energy consumption; support vector machine; metaheuristic optimization algorithms*

## I. INTRODUCTION

Residential energy consumption accounts for approximately 30% of the total energy used [1,2]. Consequently, the precise forecasting of energy consumption during the design phase and the assessment of retrofit options emerge as critical endeavors in the adventure for energy conservation and emissions reduction. The prediction of energy usage in buildings presents a formidable challenge, given its reliance on numerous factors, including building attributes, control, characteristics and maintenance, meteorological parameters, and occupants' behavior, among other sociological variables [3,4]. In response to this challenge, significant efforts from the scientific community, governmental entities, and industry stakeholders have spurred numerous research initiatives, resulting in various approaches, methodologies, and tools for estimating building energy performance. Building

energy simulation tools, notably those based on physics, such as Energy Plus [5], have gained widespread adoption for investigating and evaluating building energy efficiency. However, achieving precise simulations necessitates detailed building information, including specific space characteristics, which can be challenging to obtain [6]. Moreover, research has demonstrated substantial variations in outcomes based on the level of sophistication, both in terms of physical modeling and mathematical complexity, applied in the energy of building models [7].

In order to forecast total energy consumption or particular end uses, researchers have complemented physics-based models with a variety of statistical techniques [8]. By taking into account the characteristics of the building and its occupants and comparing the outcomes with simulations, regression-based approaches have become popular for forecasting energy consumption. Catalina et al. [9] used polynomial regression with a model displaying a maximum variance of 5% from simulated data across scenarios for estimating heating demand. Due to their ability to handle complicated interactions, more current advanced machine learning algorithms like Artificial Neural Networks ($ANN$) and Support Vector Machines ($SVM$) have been deployed. By taking into account the features of the building and its occupants and comparing the outcomes with simulations, regression-based approaches have become popular for estimating energy use. Xifara and Tsanas [10] demonstrated the superiority of Random Forest ($RF$) over regression in estimating heating load ($HL$) and cooling load ($CL$). The study employed a statistical machine learning framework to analyze how eight input variables affect $HL$ and $CLs$ in residential buildings. It systematically investigated the association strength between each input variable and the output variable star using classical and non-parametric statistical tools. Comparisons were made between classical linear regression and $RF$ for estimating $HL$ and $CL$. Simulations on 768 residential buildings demonstrated the capability to predict $HL$ and $CL$ with low mean absolute error deviations. Overall, the study supported the use of machine learning for accurate building parameter estimation in the context of energy-efficient design and operation. Li et al. [11] employed $SVM$ and $ANN$ to predict cooling demand, with $SVM$-based predictions having roughly half the errors compared to $ANN$ predictions when matched against simulation data. They established an hourly building $CL$ prediction model using $SVM$ and applied it to an

office building located in Guangzhou, China. The simulation results indicated that the *SVM* method outperformed the traditional back-propagation (*BP*) neural network model in terms of accuracy and generalization. Neto and Fiorelli [12] found similar performance between Energy Plus and an *ANN* model for energy consumption estimation. These studies collectively indicate the effectiveness of empirical methods in capturing complex achieving and relationship accuracy similar to or better than regression-based models compared to simulation results.

Instead of relying exclusively on simulated results, it is crucial to evaluate how well these sophisticated statistical models anticipate the precise energy performance of buildings in the future [13]. Significant differences between original design simulations and actual energy calculations have been found in prior studies. These differences have mostly been related to modeling assumptions, building quality, weather variations, operating practices, and occupant behavior [14]. A variety of opportunities exist to use cutting-edge approaches for examining the complicated relationship between building and occupant features and real energy performance through the analysis of large datasets as a result of the increasing availability of data regarding actual energy usage [15].

Numerous case studies have used sophisticated algorithms and previous data to predict the energy performance of buildings. For instance, Gonzalez and Zamarreno [16] presented a novel approach for short-term load prediction in buildings. The method relied on a specialized artificial neural network (ANN) that incorporated feedback from a portion of its outputs. The training of this ANN utilized a hybrid algorithm. The new system incorporated current and forecasted values of temperature, the current load, and the hour and day as inputs. The performance of this predictor underwent evaluation using real data and results from international contests. The obtained results demonstrated the high precision achieved with this system. Dong et al. [17] investigated the use of SVM for predicting building energy consumption in the tropical region, which is crucial for baseline model development and measurement and verification protocols. Four commercial buildings in Singapore were studied, employing weather data and monthly utility bills. SVM's performance, influenced by parameters C and ε, was analyzed using a radial basis function (RBF) kernel. Results indicated SVM's effectiveness, producing predictions with coefficients of variance under 3% and percentage errors within 4%. The study demonstrated SVM's feasibility and applicability in building load forecasting, offering valuable insights for accurate energy consumption predictions in tropical climates. Tso and Yau [18] conducted an empirical study comparing regression, decision tree, and *ANN* and decision tree models to beat regression techniques in forecasting power use in residential buildings. Collectively, these case studies demonstrate that machine learning algorithms offer dependable outcomes. They possess the ability to model non-linear relationships, and many are non-parametric, obviating the need for specific probability distribution assumptions [19]. It is important to remember that earlier examinations sometimes concentrated on a single structure or a small group of buildings in a particular place. Because of this, there is still a gap in the development of

generalized prediction models based on Machine Learning (*ML*) algorithms, making it difficult to use these techniques to fully examine the relevance of different architectural and occupant features for building energy efficiency [20].

Bashir and Alotaibi [21] underscored the crucial role of implementing effective building cooling and heating load prediction models for enhanced energy efficiency. In recent years, several research studies addressed challenges in determining efficient input parameters and developing accurate prediction models. Various data-driven approaches were proposed to optimize energy consumption systems and ensure indoor comfort. Despite existing reviews on prediction models, gaps remained in assessing cooling and heating load predictions. This study critically reviewed recent models, focusing on performance and accuracy. Comparative analysis revealed specific advantages for each model, yet shortcomings persisted in input parameters and implementation techniques. The review aimed to highlight and compare existing models' disadvantages in cooling and heating load predictions. Gong et al. [22] investigated the prediction of heating energy consumption in residential structures in Tianjin. They used a variety of methods, such as Support Vector Regression (SVR), Multilayer Perceptron (MLP), RF, and Light Gradient Boosted Machine (LGBM). The results showed that the LGBM model beat its competitors in a variety of assessment measures, demonstrating its potential for exact energy consumption predictions. Nebot and Mugica the [23] investigated prediction of heating and cooling loads in residential constructions utilizing fuzzy logic approaches such as fuzzy inductive reasoning (FIR) and adaptive neural fuzzy inference system (ANFIS). In their study, thirteen machine learning algorithms were compared to various fuzzy approaches, and SVR, ANFIS, and FIR performed better. Moradzadeh et al. [24] concentrated on estimating cooling and heating loads using SVR and MLP models. The MLP model had an incredible $R^2$-value of 0.9993 for heating load prediction, while the SVR model excelled with an R-value of 0.9878 for cooling load prediction, yielding outstanding results for their study. These findings illustrate the level of precision that may be achieved using machine learning algorithms. Karijadi and Chou [25] addressed the challenge of accurately predicting building energy consumption, which is crucial for effective building energy management systems. Due to the non-linear and nonstationary nature of energy consumption data, conventional prediction methods faced difficulties. The research introduced a novel hybrid approach, combining RF and Long Short-Term Memory (LSTM) based on Complete Ensemble Empirical Mode Decomposition with Adaptive Noise (CEEMDAN). The method transformed the original energy consumption data into components using CEEMDAN, where RF predicted the highest frequency component and LSTM predicted the rest. Combining the predictions yielded superior results compared to benchmark methods, as demonstrated in experiments using real-world building energy consumption data.

In the current study, inspired by prior successful results demonstrating the superior performance of *SVM* over other models, support vector regression-based models were developed to predict heating loads (*HL*) in buildings. Another advantage of this study is the utilization of numerous datasets,

including various input variables related to building geometry and glazing status, which were collected from previous literature for training predictive models. The predictive performance of a single SVR model was assessed, and in optimizing the training process, two distinct optimizers, namely the Transit Search Optimization Algorithm ($TSOA$) and the Coot optimization algorithm ($COA$), were employed. The predicted results of the three models were compared using performance metrics, including $R^2$, $RMSE$, $MAE$, $RSR$, and $MRAE$. Subsequently, the most optimal hybrid model for predicting $HL$ in buildings was determined.

The novelty of this study lies in its application of SVR models for predicting HL in buildings, driven by prior evidence showcasing the superior performance of SVR over other modeling techniques. Additionally, the study introduces a unique aspect by incorporating a diverse range of datasets that encompass various input variables associated with building geometry and glazing status. These datasets, sourced from existing literature, contribute to the comprehensive training of predictive models.

Moreover, the study distinguishes itself by evaluating the predictive performance of a singular SVR model and introducing innovation in the optimization of the training process. Specifically, the study employs two distinct optimizers, the TSOA and the COA, to enhance and fine-tune the efficiency of the SVR model training. This dual-pronged approach toward model optimization adds a novel dimension to the study, contributing to its overall uniqueness in addressing the prediction of heating loads in buildings.

In Section II, data, models, and optimizers will be introduced. In Section III, the evaluation models are developed, and the metrics used for evaluation are discussed. Finally, in Section IV, the conclusion of the study is mentioned along with the limitations and future study.

## II. Materials and Methods

### A. Data Collection

To guarantee the validity and efficacy of the approaches described in this study, the availability of reliable and substantial data is crucial. The dataset created to train the intelligent models from earlier research was utilized in this investigation. This dataset provides the crucial data needed to put the suggested strategies into practice and evaluate how well they anticipate building heating needs. Eight significant factors, including relative compactness (which represents the building's surface area-to-volume ratio), roof area, surface area, wall area, overall height, orientation, glazing area (which includes glazing, frame, and sash components), and glazing area distribution, have an impact on the analysis of the input parameters in this study. The key criteria used for the statistical analysis of the dataset are shown in Table I, together with metrics such as data averages ($Avg$), standard deviation ($St. Dev.$), minimum ($Min$), and maximum ($Max$) values.

### B. Overview of Machine Learning (ML) Methods and Optimizers

*1) Support Vector Regression (SVR):* In the early steps of pattern recognition research, support vector machines (SVM) were used to identify patterns. This approach was initially proposed by Vapnik [26] and later advocated utilizing the SVM to address issues about function approximation. The SVR methodology involves a dataset comprising $\overline{N}$ elements $\{(X_i, y_i), i = 1,2, \dots, \overline{N}\}$, where N represents the amount of training examples. The variable $X_i$ denotes the $i - th$ section of an $N-$dimensional vector, where $X_i = \{x_1, x_2, \dots, x_n\} \in R^n$, and $y_i \in R$ presents the actual value related to $X_i$.

The underlying principle of the SVR involves utilizing a ML technique to chart train data opinions, denoted as precisly $X_i$, onto a nose space that typically has $l$ dimensions.

TABLE I. The Statistical Properties of the Input Mutable of Heating

| Variables | Indicators | | | | |
|---|---|---|---|---|---|
| | Category | Min | Max | Avg | St. Dev. |
| *Relative Compactness* | Input | 0.62 | 0.98 | 0.76 | 0.11 |
| *Surface Area* (m²) | Input | 514.50 | 808.50 | 671.71 | 88.09 |
| Wall Area (m²) | Input | 245.00 | 416.50 | 318.50 | 43.63 |
| *Roof Area* (m²) | Input | 110.25 | 220.50 | 176.60 | 45.17 |
| *Overall height* (m) | Input | 3.50 | 7.00 | 5.25 | 1.75 |
| *Orientation* | Input | 2.00 | 5.00 | 3.50 | 1.12 |
| *Glazing Area* (%) | Input | 0.00 | 0.40 | 0.23 | 0.13 |
| *Glazing Area Distribution* | Input | 0.00 | 5.00 | 2.81 | 1.55 |
| *Heating* (KW) | Output | 6.01 | 43.10 | 22.31 | 10.09 |

An optimized hyperplane that precisely depicts the non-linear relationship between the output and the current input independent variables is created by carefully designing the feature space. One formal way to represent the expression for *SVR* is as shown in Eq. (1):

$$f(x) = V^T \emptyset(x) + a \tag{1}$$

Here, $a$ is the variable factor, $f(x)$ is the predicted ideals, and $V$ is the $l - dimensional$ weight factor. $\emptyset(x)$ represents the arrangement of every component $(X_i)$ into a feature space with in height dimensions. Eq. (2) represents the way the $\varepsilon -$insensitive loss function is expressed.

$$|y - f(x)|_\varepsilon = \max(0, |y - f(x)| - \varepsilon) \tag{2}$$

The residual is denoted by Eq. (3) as the disparity among the definite value, $y$, and the estimated cost, $f(x)$.

$$R(x, y) = y - f(x) \tag{3}$$

The ideal model is determined by incorporating the entire remaining element within a predefined variety of $\varepsilon$, as follows:

$$-\varepsilon \leq R(x, y) \leq \varepsilon \tag{4}$$

Eq. (4) provides the hypothesis for the complete training data. Hence, the data displays the highest disparity from the hyperplane when the remaining adheres to the condition $R(x, y) = \pm\varepsilon$. The physical distance among a specific data point $(x, y)$ and the hyperplane $R(x, y) = 0$ is calculated as $|R(x, y)|/\|V^*\|$ obtained in the following manner:

$$V^* = (1, -V^T)^T \tag{5}$$

The hypothesis of this study suggests that the most significant movement among the dataset $(x, y)$ and the hyperplane $R(x, y) = 0$ can be expressed as the adjustable $\delta$. Hence, it can be inferred that the complete train dataset meets the criteria specified in Eq. (6). The attainment of the maximum value of $\delta$ implies that the *SVR* model can demonstrate the best performance of generalization.

$$|R(x, y)| \leq \delta\|V^*\| \tag{6}$$

The highest distance is reached where the value of the $R(x, y)$ generations a predefined $\varepsilon$ value. Following this, Eq. (6) can be formulated again and expressed as Eq. (7). In order to reach the maximum value of $\delta$, it is crucial to minimize $\|V^*\|$, and as $\|V^*\|^2 = \|V\|^2 + 1$, the problem of optimization is transformed into minimizing $\|V\|$.

$$\varepsilon = \delta\|V^*\| \tag{7}$$

Despite efforts made over the train phase to minimize errors in the variety of $(-\varepsilon, \varepsilon)$, there is still a possibility that specific errors may exceed this limit. Errors occurring during training that are less than $-\varepsilon$ are denoted as $\zeta_i$, whereas training errors greater than $\varepsilon$ are represented as $\zeta^*_i$. The notations $\zeta_i$ and $\zeta^*_i$ are clarified based on following equations:

$$\zeta_i = \begin{cases} 0 & R(x_i, y_i) - \varepsilon \leq 0 \\ R(x_i, y_i) - \varepsilon & others \end{cases} \tag{8}$$

$$\zeta^*_i = \begin{cases} 0 & \varepsilon - R(x_i, y_i) \leq 0 \\ \varepsilon - R(x_i, y_i) & others \end{cases} \tag{9}$$

The primary aim of *SVR* algorithm is to identify the hyperplane that yields the optimal result though reducing the disparity among the error of training and the hyperplane. This is accomplished by utilizing the $\varepsilon$ insensitive loss purpose. Eq. (10) presents the objective function for optimizing SVR.

$$minF(W, b, \zeta_i, \zeta^*_i) = \frac{1}{2}\|W\|^2 + c\sum_{i=1}^{N}(\zeta_i + \zeta^*_i) \tag{10}$$

With the restrictions:

$$y_i - W^T\varphi(x_i) - b \leq \varepsilon + \zeta_i \quad i = 1, 2, \dots, \overline{N}$$

$$W^T\varphi(x_i) + b - y_i \leq \varepsilon + \zeta^*_i \quad i = 1, 2, \dots, \overline{N}$$

$$\zeta_i \geq 0, \zeta^*_i \geq 0 \quad i = 1, 2, \dots, \overline{N}$$

Parameter C plays an essential role in achieving a equility between minimizing training errors and an optimal separation among the hyperplane space and the data points in *SVR* involves preserving an ideal margin.

In Eq. (10), the initial part penalizes excessive weight values to maintain a flat regression function. The second part balances error margins and experience risk using the $\varepsilon$ - insensitive loss function.

Upon effectively resolving the quadratic optimization, which involves disparity constraints, the factor $W$ is derived using the guidelines explained in Eq. (1) and Eq. (11).

$$W = \sum_{i=1}^{N}(\beta^*_i - \beta_i)\varphi(x_i) \tag{11}$$

To calculate $\beta^*_i$ and $\beta_i$, it is necessary to solve a quadratic programming problem that identifies the Lagrangian multipliers.

Eq. (12) represents the SVR function:

$$f(x) = \sum_{i=1}^{N}(\beta^*_i - \beta_i)K(x_i - x) + b \tag{12}$$

The kernel function, denoted as $K(x_i - x)$, possesses the ability to nonlinearly project the train data onto a various characterized by a high$-$dimensional space with $l$ dimensions.

The Kernel function $K(x_i - x)$ possesses the proficiency to non-linearly project the train data onto a various with many dimensions (l-dimensions), rendering it well-suited for tackling issues associated with non-linear relationships. This is particularly valuable within the context of electrical forecasting. Fig. 1 demonstrates the schematic representation of the SVR's workflow.

Fig. 1.  Flowchart of SVR model.

*2) Coot Optimization Algorithm (COA):* The COA uses a metaheuristic optimization approach inspired by Coots' collective behaviors. Coots display various movements in water, like chain, random, leader-driven, and leader-adjusted, to reach food sources or specific locations. The COOT algorithm incorporates these behaviors and initiates by choosing a population using Eq. (13) [27]:

$$CootPos(i) = rand\,(1, N) \times (UB - LB) + LB \qquad (13)$$

$CootPos(i)$ is the spatial organizes of an individual coot, while $N$ represents the problem's dimensionality or the quatity of involved variables. $UB$ and $LB$ are the *upper* and *lower* limits of the exploration space, individually.

$$UB = [UB_1, UB_2, \dots, UB_N]\,, LB = [LB_1, LB_2, \dots, LB_N] \qquad (14)$$

Following the initial population setup, the positions of the coots are subsequently modified according to four patterns of movement.

*a) Random-Movement:* The position Q for this particular movement is firstly randomized:

$$Q = rand(1, N) \times (UB - LB) + LB \qquad (15)$$

Then, the position is updated to prevent getting stuck in local optima:

$$CootPos(i) = CootPos(i) + A \times R_2 \times (Q - CootPos(i)) \qquad (16)$$

The value $R_2$ is number within the range $[0, 1]$, and $A$ is defined as follows:

$$A = 1 - L \times (\frac{1}{Iter}) \qquad (17)$$

Here, $Iter$ is the maximum acceptable number of iterations, while $L$ represents the present iteration number.

*b) Chain-Movement:* To perform the chain program, one can determine the mean position of 2 coot birds utilizing the formula outlined in Eq. (18).

$$CootPos(i) = \frac{CootPos(i - 1) + CootPos(i)}{2} \qquad (18)$$

Here, $CootPos(i - 1)$ presents the position of the next coot in the sequence.

*c) Adjusting position following the leader:* In each group, a coot bird adjusts its location according to that of the leader, bringing the follower closer to the leader. The equation

given in Eq. (19) is used to determine the leader selection procedure.

$$K = 1 + (i \; MOD \; NL) \qquad (19)$$

$K$ represents the leader's index, $i$ denotes the follower coot's number, and $NL$ is the amount of group's leaders.

A coot's current location is getting updated applying Eq. (20):

$$CootPos(i) = LeaderPos(K) + 2 \times R_1 \times Cos(2R\pi) \qquad (20)$$
$$\times (LeaderPos(K) - CootPos(i))$$

$CootPos(i)$ denotes the position of the coot bird, while $LeaderPos(K)$ is the position of the chosen leader. $R_1$ represents a random number within $[0, 1]$, and $R$ is a random number within the range $[-1, 1]$.

*d) Leader-Movement*

Leader locations are updated with Eq. (21) aimed at shifting from local optimal locations to global optimal positions.

$$LeaderPos(i)$$
$$= \begin{cases} B \times B_3 \times Cos(2\pi R) \times (gBest - LeaderPos(i)) + gBest & B_4 < 0.5 \\ B \times B_3 \times Cos(2\pi R) \times (gBest - LeaderPos(i)) - gBest & B_4 \geq 0.5 \end{cases}$$
$$(21)$$

The symbol $gBest$ represents the optimal attainable position, while $B_3$ and $B_4$ are randomly chosen numbers selected from the interval $[0, 1]$. B is determined using the Eq. (22).

$$B = 2 - L \times (\frac{1}{Iter}) \qquad (22)$$

*3) Transit Search Optimization Algorithm (TSOA):* In the TSOA algorithm, there are two key parameters: the number of host stars ($ns$) and the signal-to-noise ratio ($SN$), determined based on the transit model. Noise is estimated using standard deviation from observations outside the transit phase. The product of $ns$ and SN sets the initial population size for TS [28].

This section discusses the five crucial phases of the TSOA as follows:

*a) Galaxy phase:* The algorithm begins by opting a galaxy and a random center within the search space. It then identifies habitable zones (life belts) within the galaxy by evaluating $ns$*SN random regions. $L_R$ using Eq. (23) to Eq. (25). The top $ns$ regions with the best fitness, indicating a high probability of hosting life, are selected for further algorithmic steps.

$$L_{R,I} = L_{Galaxy} + D - Noise \qquad I = 1, ..., (ns * SN) \; (23)$$

$$D = \begin{cases} c_1 L_{Galaxy} - L_r & if \; z = 1 \; (negative \; region) \\ c_1 L_{Galaxy} + L_r & if \; z = 2 \; (positive \; region) \end{cases} \quad (24)$$

$$Noise = (c_2)^3 L_r \qquad (25)$$

In the equations mentioned above, $L_{Galaxy}$ denotes the central position of the galaxy, while $L_r$ represents a randomly selected location within the exploration space. Additionally,

there are *two* coefficients, both ranging from 0 to 1, which denote a random number ($c_1$) and a random vector ($c_2$) with a dimension equal to the number of variables in the optimization. Parameter $D$ quantifies the difference between the study's context and the galaxy's center, whether in the front (positive) or back (negative) region. The zone parameter ($z$) is a random number (1 or 2) for precise positioning. To enhance accuracy, the noise parameter is applied to filter signal-related noise. A power of 3 is applied to the coefficient $c_2$ to minimize its computational impact, as noise levels are expected to be relatively close to the intended scenarios.

In the subsequent step, the algorithm selects one star from each previously identified region, corresponding to a stellar system, using Eq. (26) to Eq. (28). Consequently, at this stage, the algorithm has $ns$ stars to explore. The positions of these stars are represented as $L_s$ in Eq. (26). Notably, coefficients $c_3$ and $c_4$ in these Eqs. are random numbers ranging from 0 to 1, while the coefficient $c_5$ is a random vector with values in $[0,1]$ interval.

$$L_{S,I} = L_{R,I} + D - Noise \qquad I = 1, ..., ns \; (26)$$

$$D = \begin{cases} c_4 L_{R,I} - c_3 L_r & if \; z = 1 \; (negative \; region) \\ c_4 L_{R,I} + c_3 L_r & if \; z = 2 \; (positive \; region) \end{cases} \quad (27)$$

$$Noise = (c_5)^3 L_r \qquad (28)$$

*b) Transit phase:* In the TSOA, categorizing stars by class is essential. Therefore, the algorithm approximates the star's luminosity using Eq. (29):

$$L_I = \frac{R_I/ns}{(d_I)^2} \qquad I = 1, ..., ns \quad and \quad R_I \in \{1, ..., ns\} \qquad (29)$$

$$d_I = \sqrt{(L_s - L_T)^2} \qquad I = 1, ..., ns \qquad (30)$$

Here, $L_I$ represents star $I's$ luminosity while $R_I$ denotes its rank. $d_I$ signifies the distance among the telescope and star $I$. The telescope's location, $L_T$, is randomly selected at the outset of the algorithm and remains constant throughout optimization. To update the received light from a star, the algorithm adjusts $L_s$ by applying Eq. (31) to Eq. (33). In these equations, coefficients $c_6$ and $c_7$ are assigned random values: $c_6$ ranges from -1 to 1, and $c_7$ is a random vector with values between 0 and 1.

$$L_{S,new} = L_{s,I} + D - Noise \qquad I = 1, ..., ns \qquad (31)$$

$$D = c_6 L_{s,I} \qquad (32)$$

$$Noise = (c_7)^3 L_s \qquad (33)$$

Ultimately, the star's brightness is computed based on the newly obtained $f_s$ using the updated $L_{S,new}$. Subsequently, the new luminosity, $L_{I,new}$, is determined according to Eq. (34).

$$L_{I,new} = \frac{R_{I,new}/ns}{(d_{I,new})^2} \qquad I = 1, ..., ns \quad and \quad R_I \in \{1, ..., ns\} \qquad (34)$$

The potential for a transit event can be ascertained by comparing $L_I$ with $L_{I,new}$. The transit probability, denoted as $P_T$, is determined using Eq. (35), where it takes on values of 1 (indicating a probability of transit) or 0 (indicating no transit).

*(IJACSA) International Journal of Advanced Computer Science and Applications,*
*Vol. 15, No. 1, 2024*

If $P_T = 1$, the algorithm proceeds with the planet phase; otherwise, it executes the neighbor phase in the current iteration.

$$\begin{aligned} if \ \ L_{I,new} < L_I \qquad & P_T = 1 (Transit) \\ if \ \ L_{I,new} \geq L_I \qquad & P_T = 0 (No \ Transit) \end{aligned} \qquad (35)$$

*c) Exploitation phase:* In the Exploitation phase of the TSOA, the focus shifts to the planet's characteristics and potential for hosting life. Here, $L_P$ ($L_E$) pertains to the planet's attributes, including density, composition, and atmosphere. New knowledge ($K$) is incorporated to modify the planet's characteristics SN times (where $j = 1, \ldots, SN$) using Eq. (36) and Eq. (37). These equations involve random coefficients, such as $c_{15}$, $c_{16}$, and $c_{17}$, and parameter P, which specifies a random exponent between 1 and ($ns * SN$). Additionally, $c_K$ signifies the knowledge index with 1, 2, 3, or 4 values.

The algorithm's global solution is determined by selecting the best planet among all $ns$ detected planets.

$$L_{E,j} = \begin{cases} c_{16}L_P + c_{15}K & if c_K = 1 \ (state1) \\ c_{16}L_P - c_{15}K & if c_K = 2 \ (state2) \\ L_P - c_{15}K & if c_K = 3 \ (state3) \\ L_P + c_{15}K & if c_K = 4 \ (state4) \end{cases} \qquad (36)$$

$$K = (c_{17})^P L_r \qquad (37)$$

### III. RESULTS AND DISCUSSION

#### A. Prediction Performance Analysis

In this study, the SVR machine learning model was developed to predict HL. Furthermore, the study utilized two efficient optimization algorithms, TSOA and COA, to create hybrid SVR models, enhancing the capacity for fine-tuning model parameters. The dataset was divided into *three* subsets: train, validation, and test, with 70% of the data used for train, 15% for validation, and 15% for test [29]. The performance of these models was comprehensively assessed in Table II by comparing various metrics, including $R^2$ (coefficient of determination), $RMSE$ (Root Mean Square Error), $MAE$ (Mean Absolute Error), $RSR$ (Root Standard Ratio), and $MRAE$ (Mean Relative Absolute Error), as defined in Eq. (38) to Eq. (42) [30]:

$$R^2 = \left( \frac{\sum_{i=1}^{n}(T_i - \bar{T})(P_i - \bar{P})}{\sqrt{[\sum_{i=1}^{n}(T_i - \bar{P})^2][\sum_{i=1}^{n}(P_i - \bar{P})^2]}} \right)^2 \qquad (38)$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(P_i - T_i)^2}{n}} \qquad (39)$$

$$RSR = \frac{RMSE}{\sqrt{\frac{1}{n}\sum_{i=1}^{n}(T_i - \bar{T})^2}} \qquad (40)$$

$$MAE = \frac{1}{n}\sum_{i=1}^{n}\|P_i - T_i\| \qquad (41)$$

$$MRAE = \frac{1}{n}\sum_{i=1}^{n}\frac{|T_i - P_i|}{|T_i - \bar{T}|} \qquad (42)$$

where, $n$ is the number of samples, $P_i$ and $T_i$ are the predicted and test results, respectively. $\bar{T}$ and $\bar{P}$ are the average of the test and prediction result values.

#### B. Evaluation of Developed Models

The subsequent discussion provides a thorough examination of the model's effectiveness in predicting HL:

- The SVCO hybrid model demonstrated remarkable performance with maximum $R^2$ values of $R^2_{train} = 0.994$, $R^2_{validation} = 0.989$ and $R^2_{test} = 0.984$. These high $R^2$ values signify a strong fit between the model and the data, underscoring the reliability of the chosen input variables as robust predictors of the expected output. Also, for both hybrid models, $R^2$ for the testing stage is *lower* than that for the train stage, which indicates inadequate training performance of developed models.

- Regarding error metrics, which encompass RMSE, MAE, and MRAE, it is evident that the SVCO model *exhibits* significantly better accuracy when compared to the other models developed, demonstrating error values that are roughly half as large as those observed for the SVR single model.

TABLE II. THE RESULT OF DEVELOPED MODELS FOR SVR

| Model | Phase | Index values | | | | |
|---|---|---|---|---|---|---|
| | | RMSE | $R^2$ | MAE | RSR | MRAE |
| SVR | Train | 1.575 | 0.977 | 1.363 | 0.155 | 0.225 |
| | Validation | 1.924 | 0.967 | 1.691 | 0.195 | 4.753 |
| | Test | 1.858 | 0.966 | 1.634 | 0.186 | 0.337 |
| | All | 1.676 | 0.973 | 1.453 | 0.166 | 0.255 |
| SVCO | Train | 0.861 | 0.994 | 0.607 | 0.085 | 0.098 |
| | Validation | 1.045 | 0.989 | 0.758 | 0.106 | 0.979 |
| | Test | 1.285 | 0.984 | 0.955 | 0.129 | 0.230 |
| | All | 0.964 | 0.992 | 0.682 | 0.096 | 0.117 |
| SVTS | Train | 1.201 | 0.988 | 0.896 | 0.118 | 0.134 |
| | Validation | 1.513 | 0.978 | 1.134 | 0.154 | 5.791 |
| | Test | 1.626 | 0.975 | 1.268 | 0.163 | 0.234 |
| | All | 1.322 | 0.984 | 0.987 | 0.131 | 0.160 |

1025 | P a g e
www.ijacsa.thesai.org

A lower RSR value as a standard deviation ratio results in more accuracy of the model, so the least RSR value of 0.085 for $SVCO_{train}$ confirms its accurate prediction performance.

### C. Comparison with Published Papers

Table III shows the comparison between the presented and published papers. The comparison between the presented model and published articles focuses on key performance metrics, namely RMSE and $R^2$. The present study exhibits competitive performance with an RMSE of 0.964 and an $R^2$ of 0.992. While RMSE is higher than in some references, the $R^2$ aligns closely with high values in the literature. Variability in results across studies underscores the need for future research to explore factors influencing predictive accuracy. The discussion emphasizes the balance between predictive accuracy and generalization, acknowledging differences in dataset characteristics, model complexity, and optimization techniques. The comparative analysis contributes valuable insights for refining and advancing predictive models for building heat demand.

TABLE III. COMPARISON BETWEEN THE PRESENTED AND PUBLISHED ARTICLES

| Articles | Index values | |
|---|---|---|
| | RMSE | $R^2$ |
| Moradzadeh et al. [31] | 0.4832 | 0.9993 |
| Roy et al. [32] | 0.059 | 0.99 |
| Gong et al. [33] | 0.1929 | 0.9882 |
| Afzal et al. [2] | 1.4122 | 0.9806 |
| Present Study | 0.964 | 0.992 |

### D. Visualizing the Performance of Models

The association between observed and expected values for the three prediction models is shown in a scatter plot in Fig. 2. Additionally, the test, validation, and training datasets' $R^2$ and RMSE values for each model are supplied individually. The data points in the plot are positioned at a 45−degree angle to the horizontal axis, about 10% above and below the bold continuous line. This alignment denotes that the models perform well in terms of prediction, which leads to greater $R^2$ values.







Fig. 2. Dispersion of evolved models.

The $R^2$ value would be one in a perfect world if all data points on the observation-prediction plot were perfectly aligned with the best-fit line. This would suggest that the model has correctly and error-free estimated all of the data. The SVCO model is shown in the illustration $R^2$ values of $R^2_{train} = 0.993$, $R^2_{validation} = 0.989$, and $R^2_{test} = 0.984$ for the heating loads. These numbers outperform those of the other models, showing that the SVCO model performs better in this situation than the other hybrid models.

As shown in Fig. 3, a stacked bar plot is employed in this academic study to compare various metrics comprehensively. This visualization technique offers a clear and concise representation of the relationships among different metrics by stacking them atop one another within individual bars. Each metric is assigned a distinct color, enabling a straightforward visual assessment of their contributions to the overall outcome. Fig. 2 displays the different models' calculated RMSE, $R^2$, and MAE values. On the basis of the measurements of RMSE and MAE, a deeper look indicates that the SVCO model exhibits reduced error rates. The SVTS and SVR models have comparably decreased error rates after SVCO. Additionally, the $R^2$ measures prediction accuracy, and the SVCO model surpasses the other generated models in the research in this regard.

Fig. 3. Stacked bar plot for comparing the metrics.

Fig. 4. Error percentage of the models based on the scatter plot.



Fig. 5. The violin diagram for error percentage of proposed models.

In Fig. 4 and Fig. 5, the error percentages for the models are displayed through scatter plots and violin diagrams, with categorization into the training, validation, and test datasets. The density of data points close to zero in Fig. 3 shows how effective the strategy is. A greater concentration near zero indicates increased effectiveness. Notably, the training instances exhibit a significant preponderance of values close to zero, predominantly attributed to the SVCO method. Notably, the maximum error values for all three models are observed during the testing phase, with peak values of 32.24%, 16.11%, and 63.25% for SVR, SVCO, and SVTS, respectively. This observation underscores the beneficial fitting ability of the SVCO method.

Furthermore, Fig. 4 confirms the high accuracy of SVCO, with 25-75% of errors confined to approximately (-5 to +5). It is also evident that errors associated with SVR and SVTS are approximately three and six times higher, respectively, compared to those associated with SVCO. This further emphasizes the superior performance of the SVCO method relative to its counterparts.

## IV. CONCLUSION

In summary, this study addresses the critical imperative for precise energy consumption forecasting and the evaluation of retrofit strategies within the framework of building energy management. The complexities in predicting building energy usage, driven by multifaceted variables, including building attributes, energy systems, weather conditions, and occupant behavior, have historically posed formidable challenges. While physics-based simulations have provided valuable insights, their accuracy hinges on data comprehensiveness and modeling intricacy. In response, this research harnesses the expanding wealth of public building energy data to explore the potential of machine learning techniques, explicitly emphasizing Support Vector Regression ($SVR$) models. The research findings underscore the exceptional performance of the $SVR$ optimized with the Coot optimization algorithm (SVCO) model, consistently outperforming its counterparts by reducing prediction errors by an average of 20% to over 50% and achieving a maximum $R^2$ value of 0.992 for heating load prediction. This highlights the substantial potential of machine learning, as SVCO exemplifies, to significantly enhance the precision of energy consumption forecasts. Consequently, it empowers decision-makers in energy conservation and retrofit strategies, contributing to the overarching goals of sustainable building operations and reduced environmental impact. The study has several limitations. These include potential challenges in generalizing findings across diverse datasets and real-world scenarios due to a singular focus on SVR models. The reliance on datasets from previous literature introduces concerns about data quality, consistency, and relevance. The sensitivity of the SVR model to hyperparameters and the impact of optimization algorithms may also affect generalizability. The study's limited scope on heating loads may restrict its applicability to broader aspects of building energy performance. Future studies in this field could enhance predictive models by exploring multi-modal predictions, dynamic and adaptive models, and incorporating diverse datasets, including real-time sensor data. The inclusion of human behavior aspects, uncertainty analyses, and the application of models for guiding energy-efficient interventions in buildings are additional avenues for investigation. Furthermore, validating predictive models in real-world settings through field studies would improve practical applicability.

## REFERENCES

[1] N. Fumo, A review on the basics of building energy estimation, Renewable and Sustainable Energy Reviews 31 (2014) 53–60.

[2] S. Afzal, B.M. Ziapour, A. Shokri, H. Shakibi, B. Sobhani, Building energy consumption prediction using multilayer perceptron neural network-assisted models; comparison of different optimization algorithms, Energy (2023) 128446. https://doi.org/10.1016/j.energy. 2023.128446.

[3] S. Page, S. Krumdieck, System-level energy efficiency is the greatest barrier to development of the hydrogen economy, Energy Policy 37 (2009) 3325–3335.

[4] X.-N. Bui, H. Moayedi, A.S.A. Rashid, Developing a predictive method based on optimized M5Rules–GA predicting heating load of an energy-efficient building system, Eng Comput 36 (2020) 931–940.

[5] A. Yezioro, B. Dong, F. Leite, An applied artificial intelligence approach towards assessing building performance simulation tools, Energy Build 40 (2008) 612–620.

[6] A.-T. Nguyen, S. Reiter, P. Rigo, A review on simulation-based optimization methods applied to building performance analysis, Appl Energy 113 (2014) 1043–1058.

[7] D.B. Crawley, J.W. Hand, M. Kummert, B.T. Griffith, Contrasting the capabilities of building energy performance simulation programs, Build Environ 43 (2008) 661–673.

[8] H. Zhao, F. Magoulès, A review on the prediction of building energy consumption, Renewable and Sustainable Energy Reviews 16 (2012) 3586–3592.

[9] T. Catalina, J. Virgone, E. Blanco, Development and validation of regression models to predict monthly heating demand for residential buildings, Energy Build 40 (2008) 1825–1832.

[10] A. Tsanas, A. Xifara, Accurate quantitative estimation of energy performance of residential buildings using statistical machine learning tools, Energy Build 49 (2012) 560–567.

[11] Q. Li, Q. Meng, J. Cai, H. Yoshino, A. Mochida, Applying support vector machine to predict hourly cooling load in the building, Appl Energy 86 (2009) 2249–2256.

[12] A.H. Neto, F.A.S. Fiorelli, Comparison between detailed model simulation and artificial neural network for forecasting building energy consumption, Energy Build 40 (2008) 2169–2176.

[13] A. Gebremedhin, Optimal utilisation of heat demand in district heating system—A case study, Renewable and Sustainable Energy Reviews 30 (2014) 230–236.

[14] P. De Wilde, The gap between predicted and measured energy performance of buildings: A framework for investigation, Autom Constr 41 (2014) 40–49.

[15] F. Masoumi, S. Najjar-Ghabel, A. Safarzadeh, B. Sadaghat, Automatic calibration of the groundwater simulation model with high parameter dimensionality using sequential uncertainty fitting approach, Water Supply 20 (2020) 3487–3501. https://doi.org/10.2166/ws.2020.241.

[16] P.A. Gonzalez, J.M. Zamarreno, Prediction of hourly energy consumption in buildings based on a feedback artificial neural network, Energy Build 37 (2005) 595–601.

[17] B. Dong, C. Cao, S.E. Lee, Applying support vector machines to predict building energy consumption in tropical region, Energy Build 37 (2005) 545–553.

[18] G.K.F. Tso, K.K.W. Yau, Predicting electricity energy consumption: A comparison of regression analysis, decision tree and neural networks, Energy 32 (2007) 1761–1768.

[19] M. Hollander, D.A. Wolfe, E. Chicken, Nonparametric statistical methods, John Wiley & Sons, 2013.

[20] B.S.A.J. khiavi; B.N.E.K.A.R.T.K. hadi Sadaghat;, The Utilization of a Naïve Bayes Model for Predicting the Energy Consumption of Buildings,

Journal of Artificial Intelligence and System Modelling 01 (2023). https://doi.org/10.22034/JAISM.2023.422292.1003.

[21] M.B. Bashir, A.A. Alotaibi, Smart buildings Cooling and Heating Load Forecasting Models, IJCSNS 20 (2020) 79.

[22] M. Gong, Y. Bai, J. Qin, J. Wang, P. Yang, S. Wang, Gradient boosting machine for predicting return temperature of district heating system: A case study for residential buildings in Tianjin, Journal of Building Engineering 27 (2020) 100950.

[23] À. Nebot, F. Mugica, Energy performance forecasting of residential buildings using fuzzy approaches, Applied Sciences 10 (2020) 720.

[24] A. Moradzadeh, A. Mansour-Saatloo, B. Mohammadi-Ivatloo, A. Anvari-Moghaddam, Performance evaluation of two machine learning techniques in heating and cooling loads forecasting of residential buildings, Applied Sciences 10 (2020) 3829.

[25] I. Karijadi, S.-Y. Chou, A hybrid RF-LSTM based on CEEMDAN for improving the accuracy of building energy consumption prediction, Energy Build 259 (2022) 111908.

[26] V.N. Vapnik, The nature of statistical learning, Theory (1995).

[27] I. Naruei, F. Keynia, A new optimization method based on COOT bird natural life model, Expert Syst Appl 183 (2021) 115352.

[28] M. Mirrashid, H. Naderpour, Transit search: An optimization algorithm based on exoplanet exploration, Results in Control and Optimization 7 (2022) 100127.

[29] Q.H. Nguyen, H.-B. Ly, L.S. Ho, N. Al-Ansari, H. Van Le, V.Q. Tran, I. Prakash, B.T. Pham, Influence of Data Splitting on Performance of Machine Learning Models in Prediction of Shear Strength of Soil, Math Probl Eng 2021 (2021) 4832864. https://doi.org/10.1155/2021/4832864.

[30] A. Botchkarev, Performance metrics (error measures) in machine learning regression, forecasting and prognostics: Properties and typology, ArXiv Preprint ArXiv:1809.03006 (2018).

[31] A. Moradzadeh, A. Mansour-Saatloo, B. Mohammadi-Ivatloo, A. Anvari-Moghaddam, Performance evaluation of two machine learning techniques in heating and cooling loads forecasting of residential buildings, Applied Sciences 10 (2020) 3829.

[32] S.S. Roy, P. Samui, I. Nagtode, H. Jain, V. Shivaramakrishnan, B. Mohammadi-Ivatloo, Forecasting heating and cooling loads of buildings: A comparative performance analysis, J Ambient Intell Humaniz Comput 11 (2020) 1253–1264.

[33] M. Gong, Y. Bai, J. Qin, J. Wang, P. Yang, S. Wang, Gradient boosting machine for predicting return temperature of district heating system: A case study for residential buildings in Tianjin, Journal of Building Engineering 27 (2020) 100950.

# Application of Ant Colony Optimization Improved Clustering Algorithm in Malicious Software Identification

Yong Qian

International School of Technical Education, Sichuan College of Architectural Technology, Deyang, 618000, China

*Abstract*—Due to the increasing threat of malware to computer systems and networks, traditional malware detection and recognition technologies face difficulties and limitations. Therefore, exploring new methods to improve the accuracy and efficiency of malware identification has become an urgent need. This study introduces ant colony algorithm to optimize traditional clustering algorithms and algorithm parameters. The experimental results showed that the improvement rates of the improved algorithm in accuracy, echo value, and false alarm rate were 0.253, 0.115, and 0.056, respectively. The accuracy on the training and validation sets continued to increase and the loss curve continued to decrease. In addition, the improved algorithm had stronger modeling ability for data feature relationships and temporal information. This is of great help in improving the recognition ability of virus and worm software. The improved algorithm had a lower occupancy rate of computing resources compared to other algorithms, but it could also effectively monitor device operation. Compared with traditional methods, this method can more accurately identify malicious software and effectively identify malicious software samples from large-scale datasets. This is of great significance for protecting computer systems and network security.

*Keywords*—*Ant colony algorithm; clustering algorithm; malicious software identification; computer security; optimization algorithm*

## I. INTRODUCTION

In today's digital age, malicious software poses a huge threat to computer systems and networks. Malware refers to software programs that implant, propagate, or execute malicious behavior. Its purpose may involve stealing personal information, disrupting system functionality, malicious dissemination, etc. [1-2]. The continuous evolution and increase of malicious software make it very difficult to protect computer systems and networks from attacks. Malware detection and identification technology face significant challenges [3-4]. Traditional signature and rule-based methods are no longer adequate to cope with the increasing number of malicious software. Due to the rapid mutation and diversity of malicious software, feature extraction and classification become extremely difficult. In addition, due to the concealment and diversity of malicious software, its detection and identification require a large amount of computing resources and time [5-6]. This study utilizes Ant Colony Optimization (ACO) to improve traditional clustering algorithms and optimize their parameters. ACO is a swarm intelligence algorithm that simulates ant behavior. It simulates

the behavior of ants when searching for food and establishing pathways. ACO has been widely applied in optimization problems and has achieved significant success in solving fields such as travel salesman problems and network routing optimization [7-8]. The main contribution of the research is the proposal of an improved clustering algorithm based on ant colony optimization algorithm (ACO-CA), which is specifically designed to address the complexity of malware identification. This is the first time the ACO algorithm has been applied to the clustering problem of malicious software. Necessary adjustments and optimizations are made to the algorithm to adapt to this specific field. The ACO-CA improves the dynamics and adaptability of the clustering process by introducing ant tracking and pheromone mechanisms. This enables it to effectively handle the high-dimensional nature of malware features and the uneven distribution of samples. The algorithm proposed in this paper enhances the accuracy of malware detection and improves processing speed, demonstrating its potential for identifying malware. The deployment of the ACO-CA algorithm in practical environments is also discussed in detail, providing useful guidance for future research and application.

This study is divided into six sections. Section II provides an overview of the characteristics and threats of malware, as well as the current research status of malware identification technology. Section III proposes the hierarchical ACO-CA, which provides a detailed description of how to apply the ACO algorithm to the hierarchical aggregation clustering process. Section IV conducts empirical analysis on the performance of improved clustering algorithms and malware identification, and Section V includes discussion and analysis of experimental results. Section VI summarizes the main findings and contributions of the entire study, and explores future research directions.

## II. RELATED WORKS

The research on malware identification is an important direction in the field of computer security. With the continuous increase in the number of malicious software and the continuous development of technology, the identification of malicious software has become increasingly difficult. Therefore, researchers have been striving to improve existing malware identification techniques. Hu et al. proposed a deep sub domain adaptive network with attention mechanism. The experiment showed that the average accuracy of this method reached 97.15%, and it could quickly converge without using

a large number of target domain training datasets [9]. Mario et al. classified malicious software applications from system call traces. The method demonstrated high robustness in identifying infected applications and in code conversion and major avoidance techniques [10]. Yuan et al. proposed an anomaly detection method based on a dual head neural network. After filtering, the identified samples, the accuracy of malware detection increased by 8.62% and 13.12% respectively [11]. Sanjeev et al. proposed a novel malware detection architecture that utilizes image analysis and machine learning. Numerous experiments have shown that the methods of stacking global features and stacking local features have achieved testing accuracy of 98.34% and 98.23%, respectively. On the latest malware dataset in the real world, its testing accuracy was 92.75%, with a low false alarm rate [12]. Somayyeh et al. proposed a malware detection method based on short-term and short-term memory. Case studies have shown that the model can even detect new malware with an accuracy of over 90%. In addition, the model could detect malicious software by capturing 50 connection flows, with an AUC exceeding 99.9% [13]. Mahesh et al. proposed an adaptive red fox optimization method based on convolutional neural networks to detect whether malicious software applications are benign or malicious. Comprehensive experiments have shown that the detection accuracy of this method is 97.29% [14]. Gao et al. developed a practical system called HincTI for modeling network threat intelligence and identifying threat types. Compared with the most advanced baseline methods currently available, the proposed method could significantly improve the performance of threat type identification [15].

The basic version of ACO has matured and has been widely applied. With the deepening of research, many scholars have improved the efficiency and robustness of ACO by improving algorithm parameter settings, introducing heuristic information, and improving pheromone update strategies. Chen proposed a method for balancing enterprise resource information scheduling based on improved ACO. This method had good resource information balance scheduling ability and could effectively improve resource utilization [16]. Tan et al. established a mathematical model for spot welding path planning. The simulation analysis revealed that the ACO path was improved under six different parameters, resulting in an average path length of 10357.7509 millimeters. This is in contrast to the 10830.8394 millimeters obtained by traditional algorithms. The convergence analysis of improved ACO showed that its average number of iterations is 17. Therefore, the improved ACO had higher solution quality and faster convergence speed [17]. Wang et al. proposed a progressive

randomization approach that combines exploration and utilization with improved ACO for automatic detection of snow melting on the surface of the Antarctic ice sheet. Further validation of six automatic weather stations indicated that the proposed method had higher accuracy [18]. Wang et al. proposed an improved ant colony resource scheduling algorithm. When the number of tasks reached 200, the proposed algorithm used 17.52% less execution time and 9.58% less resources than traditional algorithms, achieving a resource allocation rate of 91.65% [19]. Saemi et al. designed a Meta heuristic algorithm based on ACO. Compared with the sequential method, this algorithm provided a solution for comprehensive problems that reduced costs by 21.64% in a reasonable amount of time. In addition, for small problems, the average difference between the solution provided by the ACO algorithm and the optimal solution (exact method) was 2.96% [20].

In summary, research on malware identification and ACO is constantly developing and advancing, providing powerful solutions for computer security and optimization issues. The use of ACO-CA has potential in malware identification, but there is a lack of relevant content in existing research, so further research is still needed. This study aims to improve ACO and achieve better results in the field of malware detection and defense.

## III. A Hierarchical Clustering Model for Malware Based on Improved Ant Colony

The initial section extracts the characteristics of malicious software and classifies its dynamic features through static analysis. Then, the extracted malware features are dimensionally reduced and centered through a vector feature matrix. In the second section, a clustering objective function is designed for ACO, and finally, pheromone update and probability transfer mechanisms are introduced to optimize ACO.

### A. Malware Feature Extraction and Processing Methods

Malicious software feature extraction and processing refers to the analysis and processing of malicious software samples, extracting feature information from them, and performing corresponding processing and encoding. The purpose is to be used in security applications such as classification, detection, and protection. The common features of malware are shown in Fig. 1.

In Fig. 1, static features include file attributes (file name, path, size), file hash value, and compilation timestamp, etc. The static analysis process is shown in Fig. 2.



Fig. 1. Classification of common malware features.

Fig. 2. Signature-based static analysis process.

In Fig. 2, the static analysis process analyzes the obtained function signature information to predict code behavior and potential defects. Dynamic characteristics refer to processes such as process behavior, system calls, registry operations, etc. Network characteristics involve network communication behavior, network traffic patterns, etc. Malicious code features include malicious code structure, invocation methods, encryption, and obfuscation. Specific behavioral characteristics refer to the encryption behavior of ransomware and the monitoring behavior of spyware. Processing malware features includes encoding, normalizing, dimensionality reduction, and other operations on the extracted features to facilitate subsequent machine learning algorithms for training and classification. The characteristics of malicious software can be represented by Hash encoding, as shown in Formula (1).

$$hash\_value = hash\_algorithm(data) \qquad (1)$$

In Formula (1), $hash\_algorithm$ is the selected Hash algorithm. $data$ is the malware feature data that needs to calculate hash encoding. Assuming there is a malware feature vector $X$ containing $n$ eigenvalues, then the Min-Max normalization formula can be used to normalize the malware features, as shown in Formula (2).

$$normalized\_value = (value - Min(X))/ \\ (Max(X) - Min(X)) \qquad (2)$$

In Formula (2), $Min(X)$ and $Max(X)$ are the minimum and maximum values of each feature, respectively. $value$ is an eigenvalue in the feature vector $X$. By applying this formula, each eigenvalue can be mapped to the range of [0, 1]. Principal Component Analysis (PCA) is used to convert high-dimensional feature vectors of malicious software into low-dimensional representations, as shown in Fig. 3.

In Fig. 3, the horizontal axis represents the selected principal components, and the vertical axis represents the original data samples. Each point represents a malware sample. The original data is projected onto the selected eigenvectors after calculating the covariance matrix to obtain eigenvalues and eigenvectors, resulting in dimensionality reduced data. Assuming there are $m$ malware samples, each with $d$ features, and the feature vector matrix is $X(m \times d)$. The feature vector matrix is centralized as shown in Formula (3).

$$X' = X - mean(X) \qquad (3)$$

In Formula (3), $X'$ is the new centralization matrix. Then to calculate the covariance matrix of the centralization matrix $X'$, as shown in Formula (4).

$$C = (1/m) * X'^{\wedge}T * X' \qquad (4)$$

In Formula (4), $X^{T}$ is the transpose matrix of $X'$. $*$ represents multiplication of matrices. $1/m$ is the normalization factor, ensuring that each element of the covariance matrix is within a reasonable range. Then to calculate the eigenvalues and eigenvectors of the covariance matrix, as shown in Formula (5).

$$C * v = \lambda * v \qquad (5)$$

In Formula (5), $\lambda$ is the eigenvector corresponding to the eigenvalues and $v$. Then, the eigenvalues $\lambda$ are sorted in descending order. When sorting feature values, to use the argsort function in the NumPy library. The feature values are stored in a one-dimensional array $eig\_vals$ and the sorting index of the feature values are calculated, as shown in Formula (6).

$$sorted\_indices = np.argsort(eig\_vals) \qquad (6)$$

In Formula (6), $sorted\_indices$ is the sorted feature value index, and then the feature values are sorted according to the sorting index to obtain Formula (7).

$$sorted\_eig\_vals = eig\_vals[sorted\_indices] \qquad (7)$$



Fig. 3. Sketch of the principal component analysis of malware features.

In Formula (7), $sorted\_eig\_vals$ is the result sorted by eigenvalues. To select the eigenvector corresponding to the largest $k$ eigenvalues, usually $k < d$, as the principal component, as shown in Formula (8).

$$V = \frac{topk\_eigvecs}{\|topk\_eigvecs\|} \quad (8)$$

In Formula (8), $topk\_eigvecs$ is the first $k$ eigenvectors. The feature vector matrix $X$ is multiplied by the projection matrix $V$ to obtain the dimensionality reduced feature matrix as shown in formula (9).

$$Y = X * V \quad (9)$$

In Formula (9), the dimension of $X$ is $(m \times d)$. The dimension of $V$ is $(d \times k)$. The dimension of $Y$ is $(n \times k)$. The reduced feature matrix can be used for subsequent malware feature analysis tasks, such as malware classification and anomaly detection. The study examined the effect of modifying a single parameter on the objective function while holding all other parameters constant. To systematically identify the optimal parameter combination and enhance the accuracy and efficiency of the algorithm, the grid search method was employed.

### B. Improved Ant Colony and Hierarchical Clustering Algorithm Optimization Design

After extracting features from malicious software through feature extraction methods, ACO is used to optimize feature selection and classifier training. Finally, hierarchical clustering algorithm is used to cluster the extracted features and identify different malicious software families. The specific process is shown in Fig. 4.

Fig. 4 first extracts and processes the features of malicious software, then identifies the features of all software, and then determines whether all software has been recognized. The above steps are repeated until the ant completes the search task. Finally, to cluster the extracted features and determine the division of malware families based on the clustering results and thresholds. ACO is a heuristic optimization algorithm that simulates the foraging behavior of ant populations. Each ant selects the next location to move based on the current pheromone and heuristic information, and the ant colony routing mechanism is shown in Fig. 5.

There are two paths in Fig. 5, ABECD and ABFCD. In path BFC, as ants increase and the amount of information increases, the probability of path selection increases. In path BEC, as time increases, the amount of information decreases, and the probability of path selection decreases. To design a clustering objective function based on internal indicators, as shown in Formula (10).

$$f = \left( \frac{\sum_{i=1}^{k} D(c_i, c)}{k} \right) \Big/ \left( \frac{\sum_{j=1}^{k} \sum_{p \in C_j} D(p, c_j)^2}{k} \right) \quad (10)$$



Fig. 4. Malware identification process.



Fig. 5. Ant Colony pathfinding mechanism.

In Formula (10), $c_i$ is the centroid of the $i$-th cluster. $c$ is the centroid of all elements. $C_j$ is the $j$-th cluster. $D$ is the proximity. $k$ is the number of clusters after clustering is completed. When clustering, not only distance information is considered, but also pheromone concentration and probability random mechanisms are introduced. The new transfer probability is Formula (11).

$$p_{ij} = \frac{u_{ij}^{\alpha} \times \left(\frac{1}{d\left(C_i, C_j\right)}\right)\beta}{\sum_{l \in Node\_Left}\left(u_{il}\right)^{\alpha}\left(\frac{1}{d\left(C_i, C_j\right)}\right)\beta} \tag{11}$$

In Formula (11), $\alpha$ is the pheromone concentration factor. $\beta$ is the heuristic information factor. $u_{ij}$ is the average pheromone concentration between cluster $i$ and cluster $j$, as shown in Formula (12).

$$u_{ij} = \frac{\sum_{k \in C_i}\sum_{n \in C_j} \tau_{kn}}{m_i \times m_j} \tag{12}$$

In Formula (12), $k$ and $n$ are two points in the path. $m_i$ and $m_j$ are the number of ants in two clusters. $\tau_{kn}$ is the concentration of pheromones between points $k$ and $n$. The algorithm process is Fig. 6.

Fig. 6 considers each point as a cluster and randomly assigns a certain number of ants to these clusters based on the proximity matrix and pheromone concentration matrix. Next, using roulette wheel to select the clusters to merge and calculate the merged result based on the objective function. If the number of iterations is not less than the set number, to output the optimal clustering merge scheme and clustering tree graph. To accelerate the convergence speed of the algorithm, the pheromone matrix is modified, as shown in Formula (13).

$$\tau_{ij}^{*} = \sigma_{ij}\tau_0 \tag{13}$$

In Formula (13), $\tau_{ij}^{*}$ is the initial pheromone concentration from element $i$ to element $j$. $\sigma_{ij}$ is the concentration coefficient. $\tau_0$ is the basic pheromone concentration. After optimizing the pheromone update mechanism, it is shown in Formula (14).

$$\tau_{ij} = \left(1-\rho\right)\tau_{ij} + \varepsilon\sum_{k=1}^{m}\Delta\tau_{ij}^{k}\left(t\right) + \Delta\tau_{ij}^{*}\left(t\right) \tag{14}$$

In Formula (14), $\varepsilon$ is the coefficient of weakness. $\Delta\tau_{ij}^{*}\left(t\right)$ is an additional pheromone on an excellent path, and the calculation Formula is (15).

$$\Delta\tau_{ij}\left(t\right) = \frac{Q}{l_{ij}} \tag{15}$$



Fig. 6. Improvement of ant colony algorithm calculation flow.

In Formula (15), $l_{ij}$ is the distance between two points. $Q$ is the amount of pheromones. To simplify the parameter optimization process, the algorithm uses a probabilistic model to predict parameter performance. New parameters are then selected for testing in regions where predicted performance is improved. This approach helps balance the relationship between exploration and exploitation, leading to more efficient identification of optimal solutions.

## IV. EMPIRICAL ANALYSIS OF ACO-CA PERFORMANCE AND MALICIOUS SOFTWARE IDENTIFICATION

When testing the performance of the improved ACO algorithm, the accuracy, echo value, and false alarm rate of different algorithms are first compared. Next, the accuracy and loss curves of the improved algorithm in the test and training sets are compared, and finally, the specificity of the improved algorithm and the traditional algorithm is compared. Empirical analysis compares the recognition of malicious software and CPU resource usage using different algorithms, and finally uses the proposed algorithm to monitor a certain device.

### A. Improved Algorithm Heating Performance Test Experiment

The analysis of improving algorithm performance selected the CICAlDroid2020 dataset as the test set and CICAlDroid2019 as the training set. CICMalDroid2019 and CICMalDroid2020 are datasets used to analyze malicious Android applications, containing 5000 and 10000 Android application samples, respectively. These samples are divided into two categories: normal applications and malicious applications. Table I shows the experimental environment for this experiment.

Experiments are conducted using traditional clustering algorithms and ACO-CA on CICCalDroid2020, and their performance in indicators such as malware recognition accuracy, recall rate, and false alarm rate were compared. The results are shown in Fig. 7.

From the data in Fig. 7, ACO-CA outperforms traditional clustering algorithms in terms of accuracy, echo value, and false alarm rate. The accuracy of ACO-CA is 0.984, while the traditional clustering algorithm is 0.731. ACO-CA can more accurately classify samples and predict the categories of malware and normal software. The echo value of ACO-CA is 0.789, while the traditional clustering algorithm is 0.674. This indicates that ACO-CA performs better in terms of echo value. The false positive rate of ACO-CA is 0.345, while the traditional clustering algorithm is 0.401. The accuracy and loss curves of the improved algorithm in the test and training sets are shown in Fig. 8.

TABLE I. EXPERIMENTAL HARDWARE AND SOFTWARE ENVIRONMENT

| Typology | Configurations |
|---|---|
| Operating system | Ubuntu 22.04.1 |
| CPU Model | Intel Core i5 1240P |
| Random access memory (RAM) | 16G |
| Hard disk | 512G |
| Python | 3.9.13 |
| Python Toolkit | Numpy+Pandas+Matplotlib+Scikit-learn etc |
| Programming Environment | Vscode+jupyter |



Fig. 7. Metrics performance of different clustering algorithms.



(a) Improvement of the ACO-TCA accuracy curve    (b) Improvement of the ACO-TCA loss curve

Fig. 8. Improvement of ACO-TCA accuracy versus loss curves.

Fig. 9.   Different algorithms specific rate fitting curves.

From Fig. 8 (a), the accuracy curves of the improved algorithm show an upward trend on both the training and validation sets. This indicates that as the training progresses, the algorithm gradually improves its accuracy during the learning process. Especially when the epoch is around 80, the accuracy tends to stabilize and remains at a high level. This indicates that the improved algorithm can achieve good classification results on both the training and validation sets, and has high accuracy. Secondly, from Fig. 8 (b), the loss curves of the improved algorithm on both the training and validation sets show a downward trend, and have already entered the range below 0.05 when the epoch is 20. Although there were some fluctuations afterwards, the overall level remained relatively low, not exceeding 0.1. This indicates that the improved algorithm can effectively reduce the loss rate during the training process, and a low loss rate indicates that the model can accurately predict the category of samples. The specificity between the improved algorithm and the traditional algorithm is shown in Fig. 9.

In Fig. 9 (a), the specificity decrease rate of traditional clustering algorithms is 0.757, while the specificity decrease rate of ACO-CA is 0.108. ACO-CA performs better at the rate of decrease in specificity and decreases more slowly. This indicates that ACO-CA can maintain a higher level of specificity when processing more samples. In Fig. 9 (b), the longitudinal intercept of traditional clustering algorithms is 98.654, while the longitudinal intercept of ACO-CA is 99.124. Therefore, ACO-CA has higher specificity values when the sample size is small.

### B. Empirical Experiment on Malicious Software Identification and Classification

Malware can be classified into various types, including viruses, worms, trojans, spyware, and adware, etc. This malicious software may steal users' personal information, damage system files, and indiscriminately send spam, posing serious risks and losses to computer systems and users. The goal of malware identification and classification is to accurately identify unknown software samples as malware or legitimate software by constructing an accurate classification model. The identification of malicious software in a dataset using different algorithms is Fig. 10.



Fig. 10.  Different algorithms for malware recognition statistics.

In Fig. 10 (a), it can be concluded that among traditional clustering algorithms, Trojan software has the highest recognition rate, accounting for 38% of the total, followed by advertising software and worm software. The recognition rates of viruses and spyware are relatively low, accounting for 13% and 8% respectively. In Fig. 10 (b), by using ACO-CA, the recognition rates of viruses and worm software have been improved, accounting for 22% and 28% respectively. The recognition rate of Trojan software has decreased, accounting for only 21%. The recognition rate of advertising software and spyware remains unchanged. This indicates that ACO-CA can improve the recognition ability of viruses and worm software in certain aspects. In Fig. 10 (c), Compared with traditional clustering algorithms, the CNN algorithm has a recognition rate of 39% for viruses and 16% for advertising software, which is significantly higher than other algorithms. In Fig. 10 (d), in the identification of malicious software based on RNN algorithm, the recognition rates of viruses and advertising software are relatively high, accounting for 27% and 26% respectively. The CPU resources occupied by different algorithms for recognition are shown in Fig. 11.

According to Fig. 11, based on ACO-CA and traditional clustering algorithms, the CPU resource utilization remains at a relatively low level (0.20-0.22 and 0.22-0.20) during the first 12 minutes of runtime. This may be the stage of algorithm initialization and data preprocessing, requiring less computational resources. After a running time of 12 minutes, the CPU resource utilization based on ACO-CA and traditional clustering algorithms rapidly increases, reaching levels of 0.78 and 0.75, respectively. The CPU resource utilization rate based on CNN algorithm and RNN algorithm is relatively high in the first 12 minutes of runtime (0.24 and 0.18). This may be because these two algorithms typically require more computing resources for convolution and loop operations. In the following period, the CPU resource utilization rate based on CNN algorithm and RNN algorithm gradually increases and stabilizes at a higher level (0.80). The real-time monitoring of the software operation of a certain device by ACO-CA is shown in Fig. 12.



Fig. 11. Variation curve of CPU resources occupied by different algorithms.



Fig. 12. Detection diagram of a device compromised by malware.

According to the data in Fig. 12, it is concluded that worm like malware has the highest number of attacks on this device. When the running time reaches 600 minutes, the number of attacks reaches 46. This may be because the security measures of the device are weak, allowing worm like malware to easily invade and attack. Trojan viruses are the second most common type of malware, with attacks fluctuating between 20 and 30 times. This may mean that the device has some protection measures in terms of security, but it is still vulnerable to Trojan virus attacks. Viruses are also a type of malware that attacks more frequently when running for 600 minutes, reaching 27 times. Advertising software and spyware have relatively fewer attacks during this period, with an average of less than 10 attacks. This may be because the purpose of advertising software and spyware is different. The former mainly obtains revenue through pop-up ads and other methods, while the latter is mainly used to monitor user activities and steal confidential information. The device may have taken certain measures to suppress the intrusion of adware and spyware.

## V. RESULTS AND DISCUSSION

Malicious software identification plays a crucial role in the field of network security, and traditional clustering algorithms have certain limitations in dealing with complex and ever-changing malicious software. Therefore, the study proposed to use the ACO algorithm to improve clustering methods, aiming to improve the accuracy and efficiency of malware detection. The algorithm has been improved to optimize the clustering process by simulating the behavioral characteristics of ant colonies. The results showed that the improved algorithm achieved a clustering accuracy of 0.984, far exceeding the traditional algorithm's 0.731. The false alarm rate has also been reduced from 0.401 in traditional algorithms to 0.345, indicating a significant improvement in reducing false positives. The accuracy improvement stabilized gradually during the training process, and the loss rate dropped below 0.05 by epoch 20. The specificity of the improved ACO algorithm decreased to only 0.108, compared to the traditional algorithm's 0.757. When processing small sample data, its specificity also showed a high level. In contrast, although CNN algorithm performed well in identifying specific categories of malware and RNN algorithm

also demonstrated good classification ability, improved ACO algorithm was more outstanding in overall performance. From a resource consumption perspective, the improved ACO algorithm maintained a low CPU usage rate during the initial 12 minutes. Although it increased thereafter, the trend of higher CPU usage in the early stages and stable increase in the later stages is relatively mild as compared to the CNN and RNN algorithms. However, there are still shortcomings in the research. Although the improved ACO algorithm has shown advantages in multiple indicators, its ability to recognize different types of malware needs further improvement. Future work will focus on optimizing algorithms for universality and efficiency on large-scale datasets, providing stronger technical support for the dynamic identification of malicious software.

## VI. Conclusion and Future Work

The constant increase of malware poses a serious threat to computer systems and network security. To solve this problem, the improved ACO algorithm was used to enhance the traditional clustering algorithm, aiming at improving the performance in the task of malware recognition. The study was processed in two stages. Firstly, the malware feature was extracted by static analysis, and the dynamic feature set of the malware was collected. Then the obtained features were processed by dimensionality reduction and centralized by vector eigenmatrix. Secondly, the clustering objective function suitable for ACO algorithm was designed, introducing pheromone updating and probability transfer mechanism to optimize the clustering effect. The experimental results showed that compared with the traditional clustering algorithm, the improved ACO-CA has achieved better performance in three aspects of accuracy, echo value and false positive rate. The improvement rates were 0.253, 0.115 and 0.056, respectively. The improved algorithm achieved recognition rates of 22% and 28% for virus and worm software, respectively. Compared with the traditional algorithm, the specificity decline rate was only 0.108, indicating a slow decline trend in maintaining specificity. From this, this paper had potential application value in improving clustering performance and reducing false alarms, providing a feasible technical approach for real-time monitoring of device software operation status. The proposed algorithm can be deployed in network security intrusion detection systems to monitor and identify malicious software behavior in real-time. The real-time monitoring capability of algorithms can be integrated into existing firewalls, intrusion detection systems, and intrusion defense systems. However, there are still shortcomings. Further research directions can combine ACO algorithm with other machine learning algorithms to further enhance the accuracy and generalization ability of malware identification, providing more specific deployments for practical applications.

## References

[1] Nani L.Y.F, Aziah A, Masnida H. A Dynamic Malware Detection in Cloud Platform. International Journal of Difference Equations, 2020, 15(2): 243-258.

[2] Ahmad A. Packing resistant solution to group malware binaries. International Journal of Security and Networks, 2020, 15(3): 123-132.

[3] Amanul I, Fazidah O, Nazmus S, Hafiz M.H.B. Prevention of Shoulder-Surfing Attack Using Shifting Condition with the Digraph Substitution Rules. Journal of Computational and Cognitive Engineering, 2023, 1(1): 58-68.

[4] Aslan, O, Samet R. A Comprehensive Review on Malware Detection Approaches. IEEE Access, 2020, 8(1): 6249-6271.

[5] Saleh M.S, Ahmed H.E.F, Mohamed S.T, Tamer H.F, Nesrin A.A. Android Malware Prevention on Permission Based. International Journal of Applied Engineering Research, 2020, 15(1): 5-11.

[6] Mouhamed B.B, Yanjun Q, Clement K.K, Kevin M.N. Evaluation of Factors Affecting Road Maintenance in Kenyan Counties Using the Ordinal Priority Approach. Journal of Computational and Cognitive Engineering, 2023, 2(3): 260-265.

[7] Monojit D, Arnab D, Avishek B, Ujjwal K.K, Samiran C. Construction of Efficient Wireless Sensor Networks for Energy Minimization Using a Modified ACO Algorithm. International Journal of Sensors, Wireless Communication and Control, 2021, 11(9): 928-950.

[8] Kumar D, Jha V.K. An improved query optimization process in big data using ACO-GA algorithm and HDFS map reduce technique. Distributed and parallel databases, 2021, 39(1): 79-96.

[9] Hu X, Zhu C, Cheng G, Li R, Wu H, Gong J. A Deep Subdomain Adaptation Network with Attention Mechanism for Malware Variant Traffic Identification at an IoT Edge Gateway. IEEE internet of things journal, 2023, 10(5): 3814-3826.

[10] Mario L.B, Marta C, Fabrizio M.M. Data-aware process discovery for malware detection: an empirical study. Machine learning, 2023, 112(4): 1171-1199.

[11] Yuan C, Cai J, Tian D, Ma R, Jia X, Liu W. Towards time evolved malware identification using two-head neural network. Journal of information security and applications, 2022, 65(Mar): 1-11.

[12] Kumar S, Janet B, Neelakantan S. Identification of malware families using stacking of textural features and machine learning. Expert Systems with Application, 2022, 208(Dec): 1-18.

[13] Somayyeh F, Amir Jalaly B. Android malware detection using network traffic based on sequential deep learning models. Software: Practice and experience, 2022, 52(9): 1987-2004.

[14] Mahesh P.C.S, Hemalatha S. An Efficient Android Malware Detection Using Adaptive Red Fox Optimization Based CNN. Wireless personal communications: An Internaional Journal, 2022, 126(1): 679-700.

[15] Gao Y, Li X, Peng H, Fang B, Yu PS. HinCTI: A Cyber Threat Intelligence Modeling and Identification System Based on Heterogeneous Information Network. IEEE Transactions on Knowledge and Data Engineering, 2022, 34(2): 708-722.

[16] Chen S. A Balanced Scheduling Method of Smart City Enterprise Resource Information Based on Improved Ant Colony Algorithm. Journal of Testing and Evaluation: A Multidisciplinary Forum for Applied Sciences and Engineering, 2023, 51(3): 1265-1276.

[17] Tan Y, Ouyang J, Zhang Z, Lao Y, Wen P. Path planning for spot welding robots based on improved ant colony algorithm. Robotica: International journal of information, education and research in robotics and artificial intelligence, 2023, 41(3): 926-938.

[18] Wang X, Guo Z, Zhang H, Wang C, Wang Y. Snowmelt detection on the Antarctic ice sheet surface based on XPGR with improved ant colony algorithm. International journal of remote sensing, 2023, 44(1/2): 142-156.

[19] Wang Y, Liu J, Tong Y, Yang Q, Liu Y, Mou H. Resource scheduling in mobile edge computing using improved ant colony algorithm for space information network. International journal of satellite communications and networking, 2023, 41(4): 331-356.

[20] Saemi S, Komijan A.R, Tavakkoli M.R, Fallah M. Solving an integrated mathematical model for crew pairing and rostering problems by an ant colony optimisation algorithm. European Journal of Industrial Engineering, 2022, 16(2): 215-240.

# The Application of MIR Technology in Higher Vocational English Teaching

Xiaoting Deng*

Department of Tourism & Foreign Languages, Henan Institute of Economics and Trade, Zhengzhou 450000, China

*Abstract*—**The traditional teaching model is teacher centered, with conservative textbooks and methods. To some extent, multimedia information retrieval technology can provide relevant information based on user query conditions, thereby alleviating the problem of information overload. This study applies image retrieval, audio and video retrieval techniques from multimedia information retrieval technology to vocational English education. It is recommended to include visual, auditory, and video materials in the course plan to meet the needs of all students. This will help ensure that the teaching objectives of each unit are achieved. Multimedia information retrieval technology may create a new learning mode in which vocational college students can use a series of mobile terminals for learning activities at anytime and anywhere, making learning more comfortable and personalized. A random double-blind survey questionnaire was designed to investigate student satisfaction and evaluate the effectiveness of multimedia information retrieval technology in vocational English teaching, in order to test the effectiveness of multimedia information retrieval technology in vocational college English teaching. According to the survey results, the majority of students acknowledge the performance of multimedia information retrieval technology in English teaching. Therefore, the application of multimedia information retrieval technology in vocational English teaching is conducive to cultivating students' self-learning ability and creative thinking ability. Meanwhile, multimedia information retrieval technology has improved the quality and level of information literacy education for college students.**

*Keywords—English teaching in higher vocational colleges; multimedia information retrieval technology; applied research; modern teaching models*

## I. INTRODUCTION

With the continued development of higher vocational education, the country has set greater standards and expectations for English teaching reform. The basic requirements of English teaching emphasize that the talents cultivated by higher vocational education are applied professionals in technology and production [1], [2]. Higher vocational English courses aim to develop students' language skills, particularly students' capacity to utilize English to deal with every day and international business activities, in addition to helping students build a strong linguistic foundation. Therefore, the ultimate goal of English learning is to cultivate learners' comprehensive English application abilities, especially their listening and speaking abilities, and to cultivate language learning as a skill [3], [4].

In recent years, higher vocational education colleges have been constantly reforming and exploring English teaching,

developing new curriculum standards and new teaching methods [5]. By integrating multimedia information retrieval (MIR) technology into English teaching, higher vocational education colleges have enhanced the teaching process and improved the effectiveness of instruction. Teachers use MIR technology to instruct English in the teaching technique. In an English lesson, the learning approach centered on grammar and reading is replaced with that based on listening and speaking [6], [7]. The teaching mode has shifted from the traditional classroom with teachers as the main body to individualized and independent learning. The classroom evaluation method has changed from the original evaluation based on grammar learning and reading to the formative evaluation based on listening and speaking ability [8].

MIR technology refers to the methods, strategies, and means of using modern information retrieval systems to retrieve relevant information. The most widely used modern information retrieval systems include online and network databases [9]. MIR technology can provide relevant information based on the user's query conditions, thereby alleviating the problem of information overload to some extent. An interesting study is applying MIR technology to English education in higher vocational colleges (HVCs). In information technology, information retrieval technology plays an important part in the transformation of the English teaching style in higher vocational colleges [10].

Colleges use network information in higher vocational and technical education to drive educational modernization. The function of MIR technology in educational development mainly includes the following aspects [11]. For students, MIR technology can produce a new way of learning, allowing students to use a variety of mobile terminals for learning activities, making learning more convenient and personalized. For colleges, MIR technology can change the mode of educational resource allocation. The relationship between users and resources is one-to-one correspondence in traditional educational resources. Still, the relationship between users and resources is transformed into a non-contact one-to-many relationship in the modern information technology environment. Using the Internet, students can share any teaching resource library to realize the fair allocation of educational information resources and promote the dynamic development and utilization of educational resources [12].

In the Internet environment, resource users are constantly adding new resources. Resources' content and expression form are constantly enriched, realizing the dynamic development of network resources. The resource user has changed from a single-user role to a dual role of user and builder. Therefore,

knowledge is updated more rapidly, and many new theories and knowledge that have not appeared in textbooks can be understood through the Internet. MIR technology can improve teachers' teaching and research abilities [13]. Through the Internet, teachers can watch and observe classes, master the latest teaching materials, interact with other teachers or education experts for discussion and exchange, and improve their teaching level. As for the college teaching mode, MIR technology promotes the innovation of the teaching mode. At present, the high-quality courses provided on the Internet can enable students to learn their favorite courses anytime and anywhere and master the learning progress. Still, there is a lack of interaction and communication. Based on MIR technology, the network teaching mode combines the benefits of these two modes to create a new teaching mode [14]. Because colleges use MIR technology to supplement traditional instruction, the role of instructors has shifted. Teachers can improve students' information retrieval level through MIR instruction to improve student's learning ability, which is conducive to improving the quality and degree of students' information literacy education [15], [16].

Traditional English teaching only attaches importance to the mechanical input and accumulation of English knowledge. Still, it ignores the inspiration of students' English learning process, especially the English language practice and other processes. Students learn passively, the classroom environment lacks vibrancy, and there is no emotional communication between teachers and students. Teachers must not only transfer knowledge and skills effectively and flexibly in organizing classroom teaching, as aided by MIR technologies. Simultaneously, teachers should engage the classroom environment, alter students' motivation to learn, and encourage emotional dialogue between teachers and students. Emotional classroom management between teachers and pupils is critical [17]. To regulate students' emotions and stimulate students' interest in learning English, teachers should reflect the teaching consciousness of democracy and equality in every link of English classroom teaching. English teaching is not simply the accumulation of words and grammar knowledge but the wisdom of allowing students to acquire knowledge and ability. In organizing teaching, teachers should release their inner passion according to students' psychology so they can independently and creatively find and analyze problems.

The evolution of MIR technology has been closely linked to the advancement of computers, databases, and networks. The potential for MIR technology in English instruction at HVCs is looking promising. Mastering MIR technology not only develops students' ability to obtain, filter, and thoroughly analyze information resources, but it also improves college students' creativity [18]. Teachers should focus on developing students' capacity to use knowledge, interact with information, and study problems independently, which can boost students' excitement for learning and improve their inventive thinking [19].

This research applies image retrieval, audio, and video retrieval technology in MIR technology to higher vocational English education to address the challenges indicated above in the traditional English teaching model. It is proposed to incorporate image, audio, and video retrieval content into the teaching objectives and develop unit teaching objectives that satisfy the needs of all students in the specific teaching process. MIR technology can create a new learning style for higher vocational students. MIR technology in colleges can potentially alter the paradigm of educational resource allocation. To evaluate the efficiency of MIR technology in higher education's English instruction. A questionnaire survey is created based on the randomized, double-blind approach to analyze the efficacy of MIR technology in higher vocational English instruction. The trial outcomes demonstrated the value of MIR technology for teaching English and motivating students in the classroom.

The research is divided into five sections. Section I elaborate and analyze the background of the application of audio and video retrieval technology in vocational English education. Section II discusses the multi-level retrieval model of English digital teaching information and the intelligent system structure of English digital teaching. Section III elaborates on modeling methods and the combination of MIR technology and education is the trend of teaching reform and development in vocational colleges. Teaching based on MIR technology is crucial for improving the quality of university teachers and changing professional teaching content. Section IV designed a questionnaire based on the random double-blind principle and surveyed student satisfaction to evaluate the effectiveness of multimedia information retrieval technology in English teaching and also delves into discussion. Section V summarizes the entire text. The research results indicate that students generally believe that traditional English teaching methods cannot meet the needs of modern talent cultivation. The above data indirectly indicates the gradual decline of traditional teaching methods and the necessity of integrating new technologies into English classroom teaching.

## II. RELATED WORKS

Cheng and Liu [20] examined the use of multimedia technology to help educators remove the barriers that prevent them from achieving their educational goals and producing graduates who meet society's requirements. They concentrate on how multimedia networks can effectively help business English majors in higher education institutions expand their job experience and strengthen their practical English abilities, as well as computer skills, communication skills, and overall cooperation skills. Zhang [21] discussed the effects of translation curriculum reform and the introduction of multimedia information retrieval technology on translation teaching. Translation teachers demand MIR technology teaching methods in language translation teaching. Liu [22] proposed a clustering method for hybrid ELT resources based on information retrieval. First, he analyzes the structural features of hybrid ELT resources and constructs a multidimensional feature distribution constraint of hybrid ELT resources by using a segmented linear fusion retrieval method. Secondly, based on the results of beam domain calculation for retrieving hybrid English teaching resources, he constructs a distribution structure model of hybrid English teaching resources. He completes the retrieval of English teaching resources. Finally, the distance between samples and each cluster prime in the set of English hybrid teaching resources is the distance within clusters, and the clustering objective

function completes the clustering of English hybrid teaching resources. A literature review of pedagogical approaches to teaching information retrieval was conducted by Fernández et al. [23] Jiang and Sun [24] created a hierarchical retrieval model for English digital teaching information and an intelligent system structure for English digital teaching. Furthermore, they used a non-negative matrix factorization approach to judge the structural similarity of English teaching material and built a hierarchical retrieval model, which increased teaching efficiency even further.

Through reviewing and analyzing existing literature, we found that although some achievements have been made in the research of blended English education resources, there are still some gaps and issues that need further exploration. Firstly, regarding the distribution structure of blended English education resources, existing research mainly focuses on the design and development of resources, while there is relatively little research on the distribution structure and models of resources. In addition, existing research lacks a comprehensive comparison and analysis of different types of educational resources, as well as in-depth research on the dynamic changes and influencing factors of the distribution structure of educational resources. Therefore, this study aims to fill this gap by conducting empirical research to explore the distribution structure model of blended English education resources, and analyzing its influencing factors and dynamic changes.

## III. MODELING METHODS

### A. MIR technology

The research of multimedia information retrieval emerged at the end of the last century and has gradually become a new important research area in information technology. MIR aims to effectively describe, organize, and find users' required multimedia information. The research of MIR involves computer vision, signal processing, pattern recognition, and many other disciplines, which have important theoretical significance [25], [26]. At the same time, MIR technology is a study area that closely integrates theory and practice. The ultimate goal is to solve information inflation and make it easier and more accurate for individuals to access the required multimedia resources [27].

Fig. 1 shows the general flow of a typical MIR system. Firstly, the system processes and analyzes the multimedia information in the database and builds the corresponding content representation and index. A canonical content query representation is generated when the user submits a retrieval requirement. Finally, the similarity is calculated according to the matching model, and the retrieval result set is returned [28], [29].

Image retrieval, audio retrieval, and video retrieval technologies are applied to English teaching in higher vocational colleges. The classification of MIR technology is shown in Fig. 2. The following describes the theories and related knowledge of the three technologies adopted in this paper.

The foundation of content-based picture retrieval is automatic feature extraction. In a broad sense, image features comprise high-level semantic information and low-level visual elements. However, existing computer vision and image comprehension technology cannot automatically extract semantic characteristics from images. Because neither picture object extraction nor recognition technology has achieved an optimal state, the most often employed low-level visual features remain low-level. The color feature is less dependent on the image's size and direction and has good compactness [30], [31].

The color histogram is the most widely used among all the color features. The color histogram is the statistics of the color values of all pixels in an image, which is defined as follows:

$$h(i) = \frac{n_i}{N}, i = 0,1, \dots, K \qquad (1)$$

where, n is the number of pixels whose color value is i in the image, n is the total number of pixels, and K is the range of possible color values. The resulting color histogram is a k-dimensional eigenvector. The spatial location of each color is unimportant to color histograms, which instead describe the proportion of various colors in the entire image. Therefore, they are especially suitable for describing the image content that does not need to consider the spatial location of specific objects.



Fig. 1. Flow chart of multimedia information retrieval.

Fig. 2. Classification of MIR techniques.

In many applications, the color image is transformed from RGB space to HSI space for retrieval [32]. The conversion formula from RGB to HSI is:

$$h = \begin{cases} cos^{-1} \frac{(r-g)+(r-b)}{2\sqrt{(r-g)^2+(r-b)(g-b)}} & b \le g \\ 2\pi - cos^{-1} \frac{(r-g)+(r-b)}{2\sqrt{(r-g)^2+(r-b)(g-b)}} & b > g \end{cases} \quad (2)$$

$$s = max(r, g, b) - min(r, g, b) \quad (3)$$

$$i = \frac{r+g+b}{3} \quad (4)$$

where, r, g, and b respectively represent the picture's pixel values of red, green, and blue.

Another very simple and effective color feature is the color moment. Unlike histograms, color moments can express certain color distribution information [33]. Let $p_{ij}$ be the i-th color component of the j-th pixel in the image, then the calculation on this color component is as follows:

$$u_i = \frac{1}{N}\sum_{j=1}^{N} p_{ij} \quad (5)$$

$$\sigma_i = (\frac{1}{N}\sum_{j=1}^{N}(p_{ij} - u_i)^2)^{\frac{1}{2}} \quad (6)$$

$$s_i = (\frac{1}{N}\sum_{j=1}^{N}(p_{ij} - u_i)^3)^{\frac{1}{3}} \quad (7)$$

Audio retrieval is a relatively new research field. Speech recognition is recognizing basic elements such as words and phrases in speech signals and then analyzing these language symbols to extract their implicit semantics. Audio retrieval is the processing, analyzing, and comprehending all audio signals, including speech and non-speech signals, to achieve the content retrieval required by users. Speech recognition technology can be applied to audio retrieval to meet speech-related retrieval needs [34].

Volume is the most widely used and easily calculated frame feature. The volume feature can be directly used for mute detection and audio segmentation. Volume is defined as follows:

$$v = \sqrt{\frac{1}{N}\sum_{i=0}^{N-1} s_i^2} \quad (8)$$

where, N is the total number of sampling points in the frame and $s_i$ is the value of each sampling point.

In general, voiceless signals have low energy and a high zero crossing rate. Therefore, by integrating zero-crossing rate and volume features, a part of voiceless speech can be prevented with low energy from being wrongly classified as silent, which is defined as follows:

$$z = \frac{1}{2}\sum_{i=1}^{N-1} |sign(s_i) - sign(s_{i-1})| \quad (9)$$

For the signal sequence $\{s_i\}$ obtained after sampling, people always want to use a model to simulate its generation. If the audio sequence $\{s_i\}$ is approximated by a linear model with finite parameters, these parameters can become important features to describe the sequence, called linear predictive coefficients. Under this linear prediction model, the prediction of the next sample can represent the weighted sum of the previous samples:

$$\hat{s}_n = \sum_{i=1}^{p} a_i s_{n-1} \quad (10)$$

where, $\{s_i\}$ is the linear prediction coefficient. In practice, it mainly establishes an optimal prediction model for the in-frame sampling sequence. Generally, the method of least mean square error is adopted. The in-frame prediction error is:

$$e = \sum_k (s_k - \sum_{i=1}^{p} a_i s_{k-1})^2 \quad (11)$$

The prediction error is set as the minimum value, and P parameters $\{s_i\}$ of the best prediction model are obtained, which are used as the linear prediction coefficient features of the current frame.

Video contains a large amount of information and rich connotation, including not only all the information of the still image but also the information of the target's movement in the scene and the information of the objective world changing with time. With the development of computer hardware and video processing software, video information has expanded rapidly. In the face of rapidly expanding video data, locating the necessary information becomes an urgent problem. To solve this problem, content-based video retrieval has been developed since the 1990s [35]. The core is to process and analyze video content to effectively obtain its content and make it easy to search and interact with data.

To obtain the global motion information, estimating the camera motion in the video is necessary. A video sequence is composed of a series of time-related image frames, and the difference between these image frames corresponds to the motion information in the video. Extracting global motion information is basically to find and determine the spatial position of each point in the background on different frames. Let $\{a_i\}$ be a set of parameters of global motion, (x, y) be a point on the current frame, and (u, v) be its corresponding point on the next frame; then the general global model can be expressed as follows:

$$\begin{cases} u = f_u(x, y, a_1, a_2, \dots) \\ v = f_v(x, y, a_1, a_2, \dots) \end{cases} \quad (12)$$

Images, video, and audio are the main components of multimedia information. Considering the richness of information in the various multimedia media, information is missing if only a single medium is used. Therefore, this paper integrates multiple MIR technology modules into higher vocational college English teaching.

### B. Application of MIR technology in higher vocational English

With the development of information technology, more and more college teachers are incorporating Internet resources into their daily teaching. The rich and colorful online information brings infinite resources to teaching and provides great teaching convenience. Resources such as pictures, audio, and videos related to the course can be used as materials for teaching. Teachers can make full use of Internet teaching tools. Such classes have the potential to increase teaching efficiency and quality significantly. As a result, it is critical for teachers to understand the fundamentals of MIR technology and to be able to access relevant information on the Internet swiftly. The combination of MIR technology and education is a trend in teaching reform and development in higher vocational institutions. Teaching based on MIR technology is critical for increasing the quality of college teachers and altering professional teaching content.

MIR is an effective way to cultivate students' autonomous learning in the English teaching of HVCs. Teachers can use pictures, audio, and video retrieval in MIR technology to promote knowledge updating, improve self-learning ability, and cultivate students' innovative spirit. Therefore, MIR effectively solves information overload, self-learning, and knowledge update problems. In practical teaching, teachers' goal design often generally describes the goal, lacking standards and levels. It only focuses on the knowledge and ability goals, ignoring students' ideological basis. Therefore, when MIR technology is used to assist English teaching, it is better to focus on the target design to conform to the content of pictures, audio, and video retrieval. The picture, audio, and video retrieval content should be integrated into the teaching objectives in the course content and teaching method. At the same time, image retrieval, audio, and video retrieval technology should be selected around the five aspects of skills, knowledge, emotion, strategy, and cultural awareness. Fig. 3 shows the curriculum content structure and teaching methods teachers formulate.

As the beacon guiding the implementation of teaching, the method of MIR should be formulated based on the pupils' actual condition. At the same time, when students use multimedia information retrieval, teachers should examine and reflect on the appropriateness of pictures, audio, and video retrieval through teaching feedback and teaching practice results. Teachers should combine teaching materials, curriculum standards, and teaching objects in English. At the same time, English teachers should take the curriculum standard as the key link, understand the essence of the textbook, and organize the textbook systematically. English teachers who work with the teaching materials can suitably modify and arrange the instructional material to satisfy various educational objectives and needs. Many English teachers in actual classroom settings are unaware of each student's unique characteristics. The polarization of pupils is caused by the fact that teaching objectives are frequently appropriate for some students but not all. Every student is unique, and English teachers should be aware of this. They should also understand the students' current circumstances and knowledge bases. At the same time, due to the student's actual level, English teachers retrieve appropriate pictures, audio, and video information and formulate unit teaching objectives to meet the needs of all students. Secondly, English teachers should pay attention to improving the effectiveness of English classroom teaching in vocational education, which requires clear objectives for classroom teaching activities. Therefore, it is possible to get the logical relationship that should be considered when setting curriculum standards, teaching objects, and textbooks, as shown in Fig. 4.



Fig. 3.   Course content and teaching method structure.

Fig. 4. Teachers set curriculum standards, teaching objects, and the logical relationship between teaching materials.

## IV. DISCUSSION AND ANALYSIS OF RESULTS

This experiment designed a questionnaire based on the principle of random double-blind and investigated students' satisfaction to evaluate the effectiveness of multimedia information retrieval technology in English teaching [36], [37]. In the questionnaire survey, the setting of the satisfaction question is shown in Fig. 5. The satisfaction settings are divided into five options, which represent "very dissatisfied," "dissatisfied," "unsure," "satisfied" and "very satisfied."

To ensure the professionalism and fairness of the survey objects, the data of this experiment were collected from teachers with rich experience in MIR technology and English teaching in HVCs as the questionnaire objects. Considering the differences in the teaching experience of different teachers, 50 students were asked to rate the classroom effect of the two teachers, respectively, and 100 student questionnaires were collected at the end. Fig. 6 shows the teaching satisfaction

evaluation of integrating MIR technology and English teaching. It can be seen that 80% are very satisfied, 10% are satisfied, 5% are uncertain, 5% are dissatisfied, and the number of very dissatisfied is 0. The results show that students agree with the teaching model of integrating MIR technology with English teaching, and students generally believe that English teaching under the new model will achieve better results. Classroom instruction combined with MIR provides students with many learning opportunities. Multimedia modes of expression also boost students' interest in learning by making the educational content more vivid and fascinating. Multimedia has a significant advantage in presenting real educational content while breaking down crucial and challenging themes.



Fig. 5. Specific questions of the questionnaire.



Fig. 6. Teaching satisfaction evaluation of the integration of MIR technology and English teaching.

Similarly, teachers with rich experience in traditional English teaching in HVCs are collected as subjects to participate in the study. Considering the differences in the teaching experience of different teachers, 50 students were asked to rate the classroom effects of the two teachers, respectively, and 100 student questionnaires were collected in the end. Fig. 7 shows the teaching satisfaction evaluation of the traditional English teaching model. It can be seen that 10% are very satisfied, 20% are satisfied, 5% are unsure, 50% are dissatisfied, and 15% are very dissatisfied. The data results show that students have low recognition of the traditional English teaching model, and students generally believe that the traditional English teaching model cannot meet the needs of modern talent training. The above data indirectly indicate the gradual decline of traditional teaching methods and the

necessity of integrating new technologies into English classroom teaching.

To ensure the credibility of the experimental results, the overall evaluation scores of the classroom effects for five students are randomly selected from the questionnaire. Furthermore, Fig. 8 compares the evaluation scores of English teachings driven by MIR technology and traditional teaching. The experimental results show that the evaluation scores of English teachings driven by multimedia technology are generally higher and more popular with students. Traditional English teaching generally has low scores and is unpopular with students. To sum up, English teaching classrooms driven by MIR technology are more popular among students than traditional English classrooms, which is the future trend of higher vocational teaching.



Fig. 7. Traditional English teaching mode of teaching satisfaction evaluation.



Fig. 8. Comparison of evaluation scores between multimedia-driven and traditional teaching models.

MIR technology has substantially enlarged the time and space restrictions of education and educational content resources. It has unprecedentedly expanded educational methods and means, altered traditional teaching modes, and ultimately led to a major reform of educational philosophy, educational concept theory, and even the entire educational system. As a result, it is vital to face the challenges of modern educational technology with courage, seize good development prospects, and improve MIR technology in higher vocational English teaching.

## V. CONCLUSIONS

Mastering MIR is a shortcut to accessing and utilizing new knowledge. By mastering MIR technology and obtaining information, English teachers in higher vocational colleges can make further effective use of information, accelerating teachers' improvement of the quality of teaching. In this paper, the image retrieval technology, audio retrieval technology, and video retrieval technology of MIR technology are applied to English teaching in HVCs, and the teaching objectives of units are formulated to meet the needs of all students. To test the effectiveness of MIR technology in higher vocational English teaching, a questionnaire is designed based on random double-blindness to investigate student satisfaction and evaluate the effectiveness of MIR technology in higher vocational English teaching. The experimental results show that 80% of the students agree with the performance of MIR technology in English teaching. 65% of the students are unsatisfied with the traditional English teaching model. Therefore, MIR technology can produce a new learning method to make learning more convenient and life-enhancing.

## FUNDING

## COMPETING OF INTERESTS

The authors declare no competing of interests.

## AUTHORSHIP CONTRIBUTION STATEMENT

Xiaoting Deng: Writing-Original draft preparation, Conceptualization, Supervision,

## AVAILABILITY OF DATA AND MATERIALS

On Request

## DECLARATIONS

Not applicable.

## REFERENCES

[1] S. Mykytiuk, T. Moroz, S. Mykytiuk, M. Moroz, and O. Dolgusheva, "Seamless Learning Model with Enhanced Web-Quizzing in the Higher Education Setting," iJIM, vol. 16, no. 03, p. 5, 2022.

[2] B. Dobreski, X. Zhu, L. Ridenour, and T. Yang, "Information Organization and Information Retrieval in the LIS Curriculum: An Analysis of Course Syllabi," Journal of Education for Library and Information Science, vol. 63, no. 3, pp. 335–350, 2022.

[3] Y. Zhang and J. Yang, "Research on Information Retrieval Effectiveness of University Scientific Researchers Based on Mental Model," Wirel Commun Mob Comput, vol. 2022, 2022.

[4] H. D. Nguyen, T.-V. Tran, X.-T. Pham, A. T. Huynh, V. T. Pham, and D. Nguyen, "Design intelligent educational chatbot for information retrieval based on integrated knowledge bases," IAENG Int J Comput Sci, vol. 49, no. 2, pp. 531–541, 2022.

[5] H. Jiang and L. Sun, "Design of Hierarchical Retrieval Model of Digital English Teaching Information Based on Ontology," Journal of Electrical and Computer Engineering, vol. 2022, 2022.

[6] N. T. Hoang and D. H. Le, "Vocational English teachers' challenges on shifting towards virtual classroom teaching," AsiaCALL Online Journal, vol. 12, no. 3, pp. 58–73, 2021.

[7] Y. Deng, "Research on the application analysis of four-dimensional teaching method in higher vocational English teaching based on big data analysis," in Journal of Physics Conference Series, 2021, p. 032059.

[8] M. Younesi and M. R. Khan, "English language teaching through the Internet at Post COVID-19 age in India: Views and Attitudes," International Journal of Research and Analytical Reviews, vol. 7, no. 3, pp. 870–875, 2020.

[9] L. Tan and F. Du, "Integrating entrepreneurship and innovation education into higher vocational education teaching methods based on big data analysis," Wirel Commun Mob Comput, vol. 2022, 2022.

[10] Abrahim, S., Mir, B. A., Suhara, H., Mohamed, F. A., & Sato, M. Structural equation modeling and confirmatory factor analysis of social media use and education. International Journal of Educational Technology in Higher Education, 16(1), 1-25, 2019.

[11] Y. Zhong, "Analysis of higher vocational English teaching behavior integrating network information teaching," in Journal of Physics: Conference Series, IOP Publishing, 2020, p. 032150.

[12] Mir, M. M., Mir, G. M., Raina, N. T., Mir, S. M., Mir, S. M., Miskeen, E., ... & Alamri, M. M. S. Application of artificial intelligence in medical education: current scenario and future perspectives. Journal of advances in medical education & professionalism, 11(3), 133, 2023.

[13] Mir, S. A., & Shakeel, D. The impact of information and communication technologies (ICTs) on academic performance of medical students: an exploratory study. International Journal of Research in Medical Sciences, 7(3), 2019, 904-908.

[14] A. Ali, "Review of Semantic Importance and Role of using Ontologies in Web Information Retrieval Techniques," International Journal of Computer and Information Technology (2279-0764), vol. 11, no. 1, 2022.

[15] L. Zhang, "The construction of digital multimedia image information retrieval model based on visual communication," in 2021 2nd International Conference on Artificial Intelligence and Information Systems, 2021, pp. 1–7.

[16] Abramova, O. V., & Korotaeva, I. E. The practical importance of student conferences in a foreign language (from the experience of working with aerospace students). Revista Espacios, 40(31), 2019.

[17] X. Feng and Y. Zhou, "English audio language retrieval based on adaptive speech-adjusting algorithm," Complexity, vol. 2021, pp. 1–12, 2021.

[18] Markova, E. M., Kuznetsova, G. V., Kozlova, O. V., Korbozerova, N. M., & Domnich, O. V. Features of the development of linguistic and communication competences of future foreign language teachers. Linguistics and Culture Review, 5(S2), 36-57, 2021.

[19] Mir, A. A., & Waheed, A. Experiences of students with disabilities in indian Higher Education: an interpretative phenomenological study, Higher Education for the Future, 9(2), 186-202, 2022.

[20] X. Cheng and K. Liu, "Application of multimedia networks in business English teaching in vocational college," J Healthc Eng, vol. 2021, 2021.

[21] Govorova, A. V., Suslova, I. P., & Shchelokova, S. V. Analysis of the online education market in Russia in the context of the theory of economic dominance. Mir novoi ekonomiki= The World of New Economy, 15(3), 77-84, 2021.

[22] J. Liu, "Clustering method of hybrid English teaching resources based on Information Retrieval," in 2022 14th International Conference on Measuring Technology and Mechatronics Automation (ICMTMA), IEEE, 2022, pp. 1015–1018.

[23] Blikstad-Balas, M., & Klette, K. Still a long way to go: Narrow and transmissive use of technology in the classroom. Nordic Journal of Digital Literacy, 15(1), 55-68, 2020.

[24] H. Jiang and L. Sun, "Design of Hierarchical Retrieval Model of Digital English Teaching Information Based on Ontology," Journal of Electrical and Computer Engineering, vol. 2022, 2022.

[25] J. Szulc, "The Features and Functions of an Effective Multimedia Information Retrieval System (MMIR)," in Multimedia and Network Information Systems: Proceedings of the 11th International Conference MISSI 2018 11, Springer, 2019, pp. 120–128.

[26] M.-H. Jang, S.-W. Kim, W.-K. Loh, and J.-I. Won, "On approximate k-nearest neighbor searches based on the earth mover's distance for efficient content-based multimedia information retrieval," Computer Science and Information Systems, 2019.

[27] M. Ravi, M. E. Naidu, and G. Narsimha, "Extracting Multimedia Information and Knowledge Discovery Using Web Mining: Challenges and Research Directions," Int. J. Applied Eng. Research, vol. 14, no. 12, pp. 2830–2836, 2019.

[28] S. Neji, L. J. Ben Ayed, T. Chenaina, and A. Shoeb, "A novel conceptual weighting model for semantic information retrieval," 07-Information Sci. Lett., vol. 10, 2021.

[29] K. V. Nayak, J. S. Arunalatha, G. U. Vasanthakumar, and K. R. Venugopal, "Design of Deep Convolution feature extraction for multimedia information retrieval," International Journal of Intelligent Unmanned Systems, no. ahead-of-print, 2022.

[30] S. Roy, S. Maity, and D. De, "MultiMICS: a contextual multifaceted intelligent multimedia information fusion paradigm," Innov Syst Softw Eng, pp. 1–19, 2022.

[31] T. M. Ghazal, M. K. Hasan, S. N. H. Abdallah, and K. A. Abubakkar, "Secure IoMT pattern recognition and exploitation for multimedia information processing using private blockchain and fuzzy logic," Transactions on Asian and Low-Resource Language Information Processing, 2022.

[32] D. Sujatha, M. Subramaniam, and C. R. Rene Robin, "A new design of multimedia big data retrieval enabled by deep feature learning and Adaptive Semantic Similarity Function," Multimed Syst, vol. 28, no. 3, pp. 1039–1058, 2022.

[33] X. Yan and J. Yan, "Design and Implementation of Interactive Platform for Operation and Maintenance of Multimedia Information System Based on Artificial Intelligence and Big Data," Comput Intell Neurosci, vol. 2022, 2022.

[34] R. Abdelmalek, Z. Mnasri, and F. Benzarti, "Audio signal reconstruction using phase retrieval: Implementation and evaluation," Multimed Tools Appl, vol. 81, no. 11, pp. 15919–15946, 2022.

[35] R. Sandeep and B. K. Prabin, "Application of Perceptual Video Hashing for Near-duplicate Video Retrieval," in Evolutionary Computing and Mobile Sustainable Networks: Proceedings of ICECMSN 2021, Springer, 2022, pp. 253–275.

[36] B. Van Tassell et al., "Rationale and design of interleukin-1 blockade in recently decompensated heart failure (REDHART2): a randomized, double blind, placebo controlled, single center, phase 2 study," J Transl Med, vol. 20, no. 1, p. 270, 2022.

[37] M. L. Kårhus et al., "Safety and efficacy of liraglutide versus colesevelam for the treatment of bile acid diarrhoea: a randomised, double-blind, active-comparator, non-inferiority clinical trial," Lancet Gastroenterol Hepatol, vol. 7, no. 10, pp. 922–931, 2022.

# Meta-Model Classification Based on the Naïve Bias Technique Auto-Regulated via Novel Metaheuristic Methods to Define Optimal Attributes of Student Performance

Zhen Ren[1], Mingmin He[2]

General Education Department, Sichuan University Jinjiang College, Meishan 620860, Sichuan, China

*Abstract*—**Accurately assessing and predicting student performance is critical in today's educational environment. Schools are dependent on evaluating students' skills, forecasting their grades, and providing customized instruction to improve their academic performance. Early intervention is essential for pinpointing areas in need of development. By predicting students' futures in particular subjects, data mining, a potent technique for revealing hidden patterns within large datasets, helps lower failure rates. These methods are combined in the field of educational data mining, which focuses on the analysis of data from educators and students with the aim of raising academic achievement. In this study, the Naive Bayes classification (NBC) model is given the main responsibility for predicting student performance. However, two cutting-edge optimization strategies, Alibaba and the Forty Thieves (AFT) and Leader Harris Hawk's optimization (LHHO), have been used to improve the model's accuracy. The study's findings show that the NBC+AFT model performs more accurately than the other models. Accuracy, Precision, Recall, and F1-Score all display impressive performance metrics for a superior model, with values of 0.891, 0.9, 0.89, and 0.89, respectively. These metrics outperform those of competing models, highlighting how successful this strategy is. Because of the NBC+AFT model's strong performance, educational institutions are getting closer to a time when they will be able to predict students' success more precisely and help them along the way, making everyone's academic journey more promising and brighter.**

*Keywords—Student performance; machine learning; classification; Naive Bayes Classification; Alibaba and the forty thieves; Leader Harris Hawk's Optimization*

## I. INTRODUCTION

Educational data mining is a powerful approach that utilizes data mining techniques to analyze vast amounts of data stored by educational institutions. These data repositories contain a wealth of information, encompassing personal and academic details of students and faculty, syllabi, question papers, circulars, and more. Various educational institutions, both universities and independent organizations, have adopted educational data mining strategies to enhance the academic experiences of their students and faculty [1], [2], [3]. These strategies are seamlessly integrated into their systems to align with their extensive databases. A few examples of educational data mining applications include:

*1)* Student performance prediction: One of the most critical aspects for educational institutions is assessing student performance. Previous academic records can serve as a basis for predicting student success. This analysis can unveil the relationships between students' abilities and interests and their academic achievements, enabling teachers to offer tailored support to those who need it the most.

*2)* Teacher evaluation: The effectiveness of teachers is often measured by their students' performance, feedback, and other relevant factors. Analyzing these data helps institutions enhance the quality of instruction and support their teaching staff better.

*3)* Question paper analysis: Evaluation of question papers can determine their level of difficulty, aiding institutions in standardizing scores across multiple sessions of examinations.

Predicting student performance is a complex challenge, akin to having a master key that opens doors to addressing underperformance by foreseeing a student's academic trajectory. This predictive ability empowers educators and decision-makers to intervene promptly and provide the necessary support to ensure every student's academic success [4]. Moreover, it extends beyond the classroom, offering insights into a student's final exam results by considering various variables, including quiz scores, homework completion, and project achievements. These holistic assessments provide a comprehensive picture of a student's academic proficiency [5], [6].

In the realm of education, machine learning algorithms have demonstrated remarkable versatility in tackling various challenges, such as classification, web mining, clustering, association rules, and deep learning. Researchers continuously explore advanced algorithms, such as clustering and classification, to develop highly accurate educational models due to the complexity of educational data. These models hold the potential to enhance the overall educational experience of students significantly.

The application of machine learning algorithms in education has yielded positive outcomes across various domains, including classification problems, clustering, association rules, web mining, and deep learning [7], [8], [9]. Researchers in the field are actively exploring advanced

algorithms like clustering and classification to build more precise educational models [10], [11], [12]. Notable examples include using machine learning techniques to correlate predictor factors with e-learning system usability, predicting students' grades, forecasting PISA test scores, and predicting adult learners' decisions to continue ESOL courses. While many prediction techniques, including regression, density estimation, and classification, are well-established, modern data science emphasizes trust and a comprehensive understanding of prediction models. Post-hoc interpretability approaches, like Local Interpretable Model-agnostic Explanations (LIME), are gaining popularity as they provide explanations for predictions made by trained black-box models, ensuring stakeholders can comprehend and rely on the insights from complex models.

In today's educational landscape, the development of robust machine-learning tools to assist educators in making well-informed decisions is not a luxury but a necessity. These tools reduce the risk of student failure, ultimately leading to improved educational outcomes. The primary objective of projects in this domain is to create dependable models for predicting student grades [13], [14]. Datasets encompassing a wide range of student performance-related factors, including personal information, educational background, and personal details, provide a comprehensive understanding of each student's situation [15], [16], [17]. By harnessing the power of machine learning, data mining, and predictive analytics, education is evolving. These technologies equip decision-makers with the tools and insights they need to identify and support underperforming students, resulting in enhanced educational outcomes for both students and institutions [18], [19].

Thammasiri et al. [20] developed a model for predicting low academic performance among freshmen. The combination of support vector machines with SMOTE yielded the greatest accuracy of 90.24%, solving class imbalance concerns. Ajay et al. [21] explored how the "CAT" social component predicts student achievement among Indians. They used four classifiers and discovered that the IB1 model had the best accuracy, at 82%. This characteristic defined people based on their social position, which had a direct influence on their educational performance. Edin Osmanbegovic et al. [22] developed a model that may predict student academic progress while addressing data dimensionality concerns. Although Naïve Bayes had the best accuracy (76.65%), it did not adequately solve the class imbalance issue. Dorina et al. [23] created a predicted model for student achievement utilizing a variety of categorization methods. While the MLP model was the most accurate at identifying successful students, it struggled to handle high-dimensional data and class imbalances. Carlos employed machine learning to develop a student failure prediction model, which achieved 92.7% accuracy using the ICRM classifier. However, due to differences in student characteristics, their study did not include testing at various educational levels.

This study is dedicated to the critical task of predicting student performance in G3 through an innovative machine-learning approach. The primary objective of this research is to optimize the performance of the Naive Bayes classification (NBC), a task made challenging by the acquisition of experimental data. The heart of this project lies in meticulous parameter optimization, which is key to enhancing the NBC model's effectiveness. To address this optimization challenge, employ a synergistic combination of two powerful algorithms: Alibaba and the Forty Thieves (AFT) and Leader Harris Hawk's optimization (LHHO). This harmonious integration of algorithms creates a cascade effect, resulting in a highly advantageous approach within the field of infrastructure, thereby elevating the intricacies involved in predicting student performance. By utilizing this novel approach, the aim extends beyond merely predicting student performance accurately. Seeks to enhance the overall effectiveness of the NBC model. Through meticulous parameter optimization and the utilization of cutting-edge algorithms, this research endeavours to provide valuable insights and solutions to the challenges faced in the education sector. This approach has the potential to revolutionize the way student performance is understood and supported on their academic journeys. Presents a promising avenue for improving the accuracy of student performance predictions and ultimately enhancing the quality of educational support. By taking these actions, it is aimed to make a positive impact on the academic prospects of all students, fostering a brighter and more promising educational future. Additionally, it strives to equip educators and decision-makers with the necessary resources to intervene effectively and guarantee the success of each student. To address the missing research summary or structure for the rest of the paper at the end of the 'Introduction' section, consider adding a brief paragraph that outlines the key components or sections to be covered in the upcoming sections of the paper.

In the following sections, a comprehensive analysis of the proposed hybrid method will delve into student performance prediction. Section I will present the experimental methodology, including the student performance data used for testing. In addition, Section II will provide an in-depth overview of the theoretical foundations, detailing the NBC coupled with AFT and LHHO. Results and comparisons with benchmark methods will be discussed in Section III, and Section IV will conclude with insights and implications drawn from the findings.

## II. MATERIALS AND METHODOLOGY

### A. Data Gathering

This section explores the application of a Naive Bayes classification (NBC) machine learning model to predict student performance based on various predictor variables. The table presents correlation coefficients that reveal the strength and direction of relationships between each predictor variable and the crucial student performance variable, G3, representing students' final grades. These correlation coefficients serve as valuable tools for understanding the multifaceted factors influencing student performance in an educational context. These three characteristics were chosen as model outputs (dependent variables), along with the number of absences from school. They were then split into four categories based on their grades: 0–12 = poor; 12–14 = acceptable; 14–16 = good; and 16–20 = excellent. Fig. 1 shows that students' ages exhibit a negative correlation with G3, indicating that older students

tends to attain lower final grades. This can be attributed to increased responsibilities and distractions accompanying age. Parental education levels of both mothers (Medu and Fedu) and fathers (Fedu) show positive correlations with student performance, suggesting that students with parents boasting higher education levels tend to achieve superior grades. Family support (famsup) and school support (schoolsup) do not reveal a significant correlation with G3, while schoolsup (school support) presents a slightly negative correlation of -0.082. Aspirations for higher education (higher) have a robust positive correlation of 0.182, highlighting the importance of nurturing academic ambitions among students. Internet access (0.098) is positively correlated with student performance, emphasizing the role of technology and online resources in augmenting learning and research opportunities. Study time (0.098) is positively correlated with G3, indicating that students who dedicate more time to studying are more likely to secure superior final grades. Previous failures (-0.360) display a pronounced negative correlation with G3, illustrating that students with fewer past failures perform substantially better, emphasizing the urgency of addressing academic setbacks promptly. These correlation coefficients provide educators, policymakers, and parents with valuable insights into the multifaceted factors influencing student performance, enabling them to tailor interventions and strategies to support students in achieving enhanced academic outcomes [24], [25]. The

potential of machine learning models, specifically NBC, in identifying these pivotal relationships and guiding evidence-based decision-making within the education sector is highlighted.

### B. Naive Bayes classification (NBC)

The Naive Bayes classification ($NBC$), a probabilistic type, employs Bayes' theorem and assumes robust feature independence. Its key strength lies in its straightforward design, obviating the need for intricate iterative parameter estimation techniques. Additionally, it has been noted by Das et al. [26] that the NB classifier is resilient to noise and irrelevant attributes. The $NB$ classifier is based on the following equation:

$$y = \underset{y_i=\{landslide,non-landslide\}}{\arg\max} P(y_i) \prod_{i=1}^{14} P\left(\frac{x_i}{y_i}\right) \qquad (1)$$

where, $P(y_i)$ is the prior probability of $y_i$, $P\left(\frac{x_i}{y_i}\right)$ is the posterior probability, and it can be calculated by:

$$P\left(\frac{x_i}{y_i}\right) = \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-(x_i-\mu)^2}{2\sigma^2}} \qquad (2)$$

where, $\mu$ is the mean and $\sigma$ is the standard deviation of $x_i$. The flowchart of the NBC is shown in the Fig. 2.



Fig. 1. Correlation matrix for the input and output variables.

Fig. 2. The flowchart of the NBC.

## C. Alibaba and the Forty Thieves (AFT)

The present study provides an explanation of the mathematical structure underlying the fundamental AFT algorithm, which is extensively detailed. Three separate states that are included in the framework can be looked at and defined as follows:

State one: Based on data gathered from a source, the thieves' chase after Ali Baba can be simulated using Eq. (3), which shows their relative positions [27], [28].

$$x_i^{t+1} = gbest^t + \left[\text{Td}^t(best_i^t - y_i^t)r_1 + \text{Td}^t\left(y_i^t - m_{a(i)}^t\right)r_2\right]sgn(rand - 0.5), \quad p \geq 0.5, \quad q > P_{p^t} \quad (3)$$

$x_i^{t+1}$ denotes the position of the $i - th$ thief at the next time step $(t + 1)$.

$m_{a(i)}^t$ represents the level of Marjaneh's wit used to disguise thief $i$, at time $t$.

$best_i^t$ represents the best position achieved by thief $i$ up to the current time step $(t)$.

$gbest^t$ refers to the best global position achieved by any thief up to the current time step $(t)$.

$r_1$, $r_2$, $rand$, $p$, and $q$ are randomly generated values that fall within the range of $[0,1]$.

$p \geq 0.5$ indicates either a value of 0 or 1.

$y_i^t$ depicts the position of Ali Baba concerning thief $i$, at time $t$.

$a$ is defined by using Eq. (6).

$sgn(rand - 0.5)$ take on a value of either -1 or 1.

$\text{Td}^t$ represents the tracking distance of the thieves, which is defined by Eq. (4).

$P_{p^t}$ represents the potential perceptual ability of the thieves to detect Ali Baba, as defined in Eq. (5).

$$\text{Td}^t = \tau_0 e^{-\tau_1(\frac{t}{T})^{\tau_1}} \quad (4)$$

$\tau_0$ $(\tau_0 = 1)$ serves as an initial estimate for the tracking distance.

$\tau_1$ $(\tau_1 = 2)$ is employed to regulate the balance between exploitation and exploration.

$t$ and $T$, respectively, refer to the current and maximum iteration values.

$$P_{p^t} = \lambda_0 \log(\lambda_1 (\frac{t}{T})^{\lambda_0}) \qquad (5)$$

$\lambda_0$ ($\lambda_0 = 1$) represents the final estimate of the probability that the thieves will succeed in achieving their goal following the search.

$\lambda_1$ ($\lambda_1 = 1$) stands for a constant that controls the ratio of exploration to exploitation.

$$a = [(n-1).rand(n,1)] \qquad (6)$$

The result of creating a sequence of random numbers between 0 and 1 is **U** and $n$, 1.

$$m_{a(i)}^t = \begin{cases} x_i^t & if\ f(x_i^t) \geq f(m_{a(i)}^t) \\ m_{a(i)}^t & if\ f(x_i^t) < f(m_{a(i)}^t) \end{cases} \qquad (7)$$

$f(0)$ represents the score or value of the fitness function.

State two: When the thieves realize they have been tricked, they might start venturing into unknown and unexpected areas.

$$x_i^{t+1} = T d^t [(u_j - l_j)r + l_j]; p \geq 0.5, q \leq P_{p^t} \qquad (8)$$

The boundaries of the search space for dimension j are represented by $u_j$ (the upper bound) and $l_j$ (the lower bound).

r is a random variable created within the range of [0, 1].

State three: The thieves may investigate additional search positions outside of those obtained by applying equations in order to enhance the AFT algorithm's exploration and exploitation components. The following scenario can be mathematically expressed as Eq. (9):

$$x_i^{t+1} = gbest^t - [Td^t(best_i^t - y_i^t)r_1 \\ + Td^t(y_i^t - m_{a(i)}^t)r_2]sgn(rand \\ - 0.5) \qquad (9)$$

Algorithm 1 provides an exact and succinct presentation of the iterative pseudo-code steps of the basic AFT algorithm.

---

**Algorithm 1: AFT algorithm**

---

Define and begin the control parameters.

Begin and evaluate the initial, best, and global positions of all thieves

Begin Marjane's wit level concerning all thieves

Set $t \leftarrow 1$

While ($t \leq T$) do

Update the parameter $P_{p^t}$ using Eq. (5).

for each thief, do

if ($p \geq 0.5$) then

if ($q \geq P_{p^t}$) then

Update the thieves' position using Eq. (4).

else

Update the thieves' position using Eq. (8).

end if

else

Update the thieves' position using Eq. (9).

---

end if

end for

Update the new, best, and global positions of all thieves

Update Marjane's wit plans using Eq. (7).

$t = t + 1$

end while

Return the best global solution

---

### D. Leader Harris Hawk's Optimization (LHHO)

The algorithm known as LHHO was developed using the exploratory behavior of the Harris hawk as a model. Owing to its equal chance $q$ perching strategy, the original Harris Hawks Optimisation (HHO) algorithm has a finite exploration capacity. If $q$ is greater than or equal to 0.5, then hawks will randomly choose a tall tree to perch on; if $q$ is less than 0.5, then they will base their perching decisions on the locations of other family members [29]. This is in accordance with the *HHO* algorithm. However, this limitation can be overcome by assigning a perch probability to each hawk [30].

During the exploration phase ($|E| \geq 1$) a concept called adaptive perch probability ($h_{ap}^i$) can be introduced for the *ith* hawk. This probability value is determined by the fitness value of the current hawk with a position vector $X_i$ denoted as $f(X_i)$, as well as the fitness values of the best-performing hawk with the position vector $X_{prey}$, denoted as $f(X_{prey})$, and the worst-performing hawk with the position vector $X_{worst}$, denoted as $f(X_{worst})$. By taking these factors into account, the adaptive perch probability ($h_{ap}^i$) can be formulated as:

$$h_{ap}^i = \frac{|f(X_i) - f(X_{prey})|}{|f(X_{worst}) - f(X_{prey})|}, \qquad i = 1,2,3,\dots,N \qquad (10)$$

Then, the exploration phase can be modeled as:

$$X_i(new) = \begin{cases} X_{round}(t) - d_1|X_{round}(t) - 2d_2 X_i(t)| & q > h_a^i \\ (X_{prey}(t) - X_w(t)) - d_3(LB + d_4(UB - LB)) & q < \end{cases} \qquad (11)$$

The model incorporates $X_w(t)$, which reflects the population's average position vector during the exploration phase of N hawks. In contrast, during the exploitation phase ($|E| < 1$), four offensive techniques that are similar to those used in HHO are employed by the model.

- Soft besiege ($r \geq 0.5\ and\ |E| \geq 0.5$)

$$X_i(new) = X_{prey}(t) - E|JX_{prey}(t) - X_i(t)| \qquad (12)$$

Where J is the jump strength as of $J = 2(1 - r_5)$

- Hard besiege ($r \geq 0.5\ and\ |E| < 0.5$)

$$X_i(new) = X_{prey}(t) - E|X_{prey}(t) - X_i(t)| \qquad (13)$$

- Soft besiege with progressive rapid dives ($r < 0.5\ and\ |E| \geq 0.5$)

$$X_i(new) = \begin{cases} Y_i & if\ f(Y_i) < f(X_i(t)) \\ Z_i & if\ f(Z_i) < f(X_i(t)) \end{cases} \qquad (14)$$

The Yi and Zi can be calculated using Eq. $Y_i = X_{prey}(t) - E|JX_{prey}(t) - X_i(t)|$ and Eq. $Z_i = Y_i + S \times LF(D)$, respectively.

- Hard besiege with progressive rapid dives ($r < 0.5$ and $|E| < 0.5$)

$$X_i(new) = \begin{cases} Y_i & if\ f(Y_i) < f(X_i(t)) \\ Z_i & if\ f(Z_i) < f(X_i(t)) \end{cases} \quad (15)$$

The equations used to calculate $Y_i$ and $Z_i$ are as follows: $Z_i = Y_i + S \times LF(D)$, respectively. It can be observed that the escape energy $|E|$ remains below 1 after 50% of the maximum iterations, indicating that the HHO algorithm only exploits solutions after this point. This restricted investigation increases the likelihood of discovering less-than-optimal solutions and becoming stuck in a local minimum. To help explore the end, a leader-based mutation-selection method is proposed as an addition to the HHO algorithm.

Here, to put the leader-based mutation-selection strategy into practice, first determine the position vectors of the best, second-best, and third-best hawks, denoted as $X_{best}^t$, $X_{best-1}^t$, and $X_{best-2}^t$, respectively. These position vectors are determined based on the fitness function value of the new position vector $X(new)$ among the $N$ individual hawks. The study can then define the new mutation position vector for the $i\_th$ hawk, denoted as $X_i(mut)$, as follows:

$$X_i(mut) = X_i(new) + 2 * \left(1 - \frac{t}{t_{max}}\right) \\ * (2 * rand - 1)(2 * X_{best}^t \\ - (X_{best-1}^t + X_{best-2}^t)) \\ + (2 * rand - 1)(X_{best}^t - X_i(new)) \quad (16)$$

where a rand is a random number in the range $(0, 1)$. Then, the position vector for the next generation $X_i(t + 1)$, can be obtained by the selection process described in Eq. (17). Similarly, the $X_{prey}$ is updated using Eq. (17). The flowchart of the LHHO is shown in Fig. 3.

$$X_i(t + 1) = \begin{cases} X_i(mut) & f(X_i(mut)) < f(X_i(new)) \\ X_i(new) & f(X_i(mut)) > f(X_i(new)) \end{cases} \quad (17)$$

$$X_{prey} = \begin{cases} X_i(mut) & f(X_i(mut)) < f(X_{prey}) \\ X_i(new) & f(X_i(new)) > f(X_{prey}) \end{cases} \quad (18)$$

*E. Performance Evaluation Methods*

Numerous evaluation criteria are used to assess the classifiers' performance. The most popular criterion for assessing classification accuracy is PESTEL, which gauges a classifier's efficacy by looking at the proportion of correctly predicted samples, as shown in the equation below. Two additional popular evaluation indices are precision and recall. The ratio of values with a positive class to those that are expected to be positive is known as recall. Conversely, precision, which can be defined as the following equations, is the likelihood that a positive prediction will come true. The f1-score, which is defined as follows, is a new value that can be produced by combining Precision and Recall.

| Equations contain | Equation | Assessment | |
|---|---|---|---|
| TP means that the outcome was positive and in line with this prediction. FP means it was not what was expected, which was a negative result. TN indicates that the outcome was negative, as predicted. FN means that although the study was expecting a negative result, the outcome was positive. | $Accuracy = \dfrac{TP + TN}{TP + TN + FP + FN}$ | *Higher is desirable* | (19) |
| | $Precision = \dfrac{TP}{TP + FP}$ | *Higher is desirable* | (20) |
| | $Recall = TPR = \dfrac{TP}{P} = \dfrac{TP}{TP + FN}$ | *Higher is desirable* | (21) |
| | $F1\ score\ = \dfrac{2 \times Recall\ \times Precision}{Recall + Precision}$ | *Higher is desirable* | (22) |

## III. RESULTS AND DISCUSSION

The outcomes of the models that were given are shown in Table I. Each model was assessed using a variety of index values, such as accuracy, precision, recall, and F1-score. The first model, NBC+AFT, demonstrated its ability to predict student performance with an accuracy of 0.891 accurately. Additionally, it showed a high degree of precision (0.9), indicating that it could correctly predict positive outcomes. The model demonstrated an F1-score of 0.89, signifying a balance between precision and recall, and a recall of 0.89, indicating its efficacy in identifying pertinent instances.

In comparison to the NBC+AFT model, the second model, NBC+LHH, performed marginally worse in terms of accuracy (0.881), precision (0.88), recall (0.88), and F1-score (0.88), but it was still very capable of making predictions. The accuracy of the third model, NBC, was 0.873, indicating that it was capable of producing accurate predictions. Additionally, it displayed F1-score, precision, and recall values of 0.87, demonstrating a balanced performance in terms of true positive predictions and the model's capacity to find pertinent instances. Overall, these models' results show how well optimization methods like AFT and LHHO can be used to improve efficiency. The NBC+AFT model slightly outperformed the others, demonstrating its potential for accurate student performance prediction, even though all models achieved high accuracy and showed a trade-off between precision and recall.

Fig. 3. Flowchart of LHHA.

The performance evaluation indices for the developed models NBC+AFT, NBC+LHH, and NBC are shown in this table. To give a complete picture of how well the models predict student performance, they are evaluated across a range of performance grades, from Excellent to Poor.

TABLE I. RESULT OF PRESENTED MODELS

| Model | Index values | | | |
|---|---|---|---|---|
| | *Accuracy* | *Precision* | *Recall* | *F1 _score* |
| NBC+AFT | 0.891% | 0.9% | 0.89% | 0.89% |
| NBC+LHH | 0.881% | 0.88% | 0.88% | 0.88% |
| NBC | 0.873% | 0.87% | 0.87% | 0.87% |

NBC+AFT:

- Poor: With a precision of 0.93, the model's ability to predict poor grades is impressive and suggests a strong identification of students who perform poorly. Additionally, it has a high recall of 0.97, indicating that it is capable of identifying most underperformers. The F1-score in this category is 0.95, which indicates a very balanced performance.

- Acceptable: Within the Acceptable grade, the NBC+AFT model showcases a precision of 0.84, signifying its capability to identify students with acceptable performance correctly. However, the recall is 0.74, indicating that it might miss some of these

students. The F1-score is 0.79, reflecting a reasonable balance between precision and recall.

- Good: The model consistently maintains a precision of 0.75 in the good grade category, demonstrating its dependability in identifying students who are performing well. Additionally, its recall score is 0.92, indicating that it can identify the majority of high performers. The F1-score of 0.83 indicates that recall and precision are in a healthy balance.

- Excellent: The NBC+AFT model exhibits high precision (1) for the Excellent grade, indicating a strong ability to recognize students who perform excellently. With a recall of 0.62, the model appears to account for 62% of students who perform exceptionally well. With an F1-score of 0.77, recall and precision are fairly balanced.

NBC+LHH and NBC:

- In all grade categories, the NBC+LHH model performs similarly, with F1-scores, precision, and recall matching those of the NBC+AFT model. Precision and recall for the NBC model are marginally different from those of the NBC+AFT and NBC+LHH models. It continues to perform well, nevertheless, in recognizing students in various grade levels.

These assessment indices offer insightful information about how well the developed models performed, highlighting how well they predicted student performance across a range of grade levels. The decision between NBC+AFT, NBC+LHH, and NBC may be influenced by particular educational environments as well as the intended harmony between recall and precision for various grade levels.

The line symbol plot shown in Fig. 4, as illustrated in the Table II, presents a visual representation of the measured data compared to the predictions generated by three distinct models: NBC+AFT, NBC+LHHO, and NBC. The measured values represent the actual number of students falling into each performance category, while the model predictions indicate the estimated numbers for each category.

*1) Poor performance:*

- NBC+AFT (226): The NBC+AFT model predicts that 226 students will perform poorly.

- NBC+LHHO (226): The NBC+LHHO model, closely aligned with NBC+AFT, also estimates 226 students to have poor performance.

- NBC (222): The standard NBC model predicts 222 students to fall into this category.

*2) Acceptable performance:*

- NBC+AFT (46): The NBC+AFT model predicts that 46 students will perform at an acceptable level.

- NBC+LHHO (45): The NBC+LHHO model estimates 45 students to have acceptable performance.

- NBC (45): The standard NBC model concurs with NBC+LHHO, also predicting 45 students in this category.

*3) Good performance:*

- NBC+AFT (55): The NBC+AFT model predicts that 55 students will achieve a good performance level.

- NBC+LHHO (51): The NBC+LHHO model estimates 51 students to fall into this category.

- NBC (49): The standard NBC model predicts 49 students in this group.

TABLE II.        EVALUATION INDEXES OF THE DEVELOPED MODELS' PERFORMANCE BASED ON GRADES

| Model | Grade | Index values | | |
|---|---|---|---|---|
| | | *Precision* | *Recall* | *F1-score* |
| NBC+AFT | Excellent | 1 | 0.62 | 0.77 |
| | Good | 0.75 | 0.92 | 0.83 |
| | Acceptable | 0.84 | 0.74 | 0.79 |
| | Poor | 0.93 | 0.97 | 0.95 |
| NBC+LHH | Excellent | 0.93 | 0.65 | 0.76 |
| | Good | 0.73 | 0.85 | 0.78 |
| | Acceptable | 0.83 | 0.73 | 0.78 |
| | Poor | 0.93 | 0.97 | 0.95 |
| NBC | Excellent | 0.85 | 0.72 | 0.78 |
| | Good | 0.74 | 0.82 | 0.78 |
| | Acceptable | 0.78 | 0.73 | 0.75 |
| | Poor | 0.94 | 0.95 | 0.94 |

Fig. 4.    Line-symbol plot for the classification accuracy of Meta-models.

*4) Excellent performance:*

- NBC+AFT (25): The NBC+AFT model predicts that 25 students will attain an excellent level of performance.

- NBC+LHHO (26): The NBC+LHHO model closely aligns with NBC+AFT, estimating 26 students to achieve excellence.

- NBC (29): The standard NBC model forecasts that 29 students will reach an excellent performance level.

These line symbol plot values illustrate how well the models align with the actual measured data for different performance categories. The variations in the predictions of each model offer insights into their individual capabilities and accuracy in identifying student performance levels. In this context, NBC+AFT, NBC+LHHO, and NBC exhibit similarities and differences in their predictions, highlighting the

strengths and limitations of each approach in assessing student performance.

Three confusion matrices that show how the NBC, NBC+AFT, and NBC+LHHO models relate to the observed and predicted classes are shown in Fig. 5. The observed classes are plotted on the horizontal axis, and the predicted classes are plotted on the vertical axis. Interestingly, these matrices' diagonal cells—which match the precise predictions—have higher values than their off-diagonal cells.

- NBC+AFT: Specifically, the NBC+AFT hybrid model shows an impressive capacity to predict most observation classes accurately. To provide more context, let's look at the NBC+AFT plot. Out of the 40 students in the excellent class, the NBC+AFT hybrid model correctly predicts 25 of them to be in the same excellent category. The remaining three students are incorrectly assigned to the poor class, 1 to the acceptable class, and 11 to the good class.

- NBC+LHHO: In the NBC+LHHO storyline, the impoverished class comprises 233 pupils. In this bad class, the NBC+LHHO hybrid model predicts 226 students with skill; only four students are incorrectly placed in the acceptable class, and only three students are incorrectly placed in the good class.

- NBC: On the other hand, the NBC story revolves around 60 pupils in the superior class. Of these, 49 are correctly predicted by the NBC hybrid model to be in the good category; one student is mistakenly placed in the poor class, five in the acceptable class, and five in the excellent class.

These results underscore the efficacy of the NBC+AFT hybrid model in accurately predicting student performance classes, with notably fewer misclassifications compared to the other models.

The convergence curve of hybrid models with 200 iterations is shown in Fig. 6. The accuracy parameter is represented by the vertical axis in this visualization, and the horizontal axis corresponds to the number of iterations. This graph's analysis reveals that the NBC+LHHO hybrid model, which records an accuracy of 0.76 and reaches its ideal iteration at number 126, is the model with the lowest accuracy. As an illustration of this, the green NBC+AFT hybrid model achieves the highest accuracy value of all the models, 0.79.9. This model performs better than the others in terms of accuracy, reaching its optimal iteration point at 128.



Fig. 5. Confusion matrix for each model's accuracy.



Fig. 6. Convergence curve of hybrid models.

## IV. CONCLUSION

Forecasting student performance is still an important task in today's educational environment. Educational establishments are responsible for determining the skills of their students, projecting their academic performance, and making proactive efforts to enhance their future success. Predictive model accuracy and efficacy have significantly increased as a result of utilizing machine learning techniques, particularly the Naive Bayes classification (NBC), in conjunction with sophisticated optimization algorithms like Alibaba and the Forty Thieves (AFT) and Leader Harris Hawk's optimization (LHHO). The improvements are especially noticeable when looking at important assessment metrics like F1-Score, Accuracy, Precision, and Recall. Forecasting student performance is still an important task in today's educational environment. Educational establishments are responsible for determining the skills of their students, projecting their academic performance, and making proactive efforts to enhance their future success. Predictive model accuracy and efficacy have significantly increased as a result of utilizing machine learning techniques, particularly the Naive Bayes classification (NBC), in conjunction with sophisticated optimization algorithms like Alibaba and the Forty Thieves (AFT) and Leader Harris Hawk's optimization (LHHO). The improvements are especially noticeable when looking at important assessment metrics like F1-Score, Accuracy, Precision, and Recall. In this thorough analysis, the NBC+AFT hybrid model has proven to be the best performer, consistently outperforming other models. With its outstanding performance in terms of Accuracy, Precision, Recall, and F1-Score, this model is the best option for educational institutions committed to improving the prediction of student performance. It performs exceptionally well at predicting academic grades with the least amount of incorrect categorizations, which is an essential feature for making informed decisions. The importance of sophisticated machine learning models and optimization strategies in the field of predicting student performance is highlighted by this study. In particular, the NBC+AFT hybrid model provides educational institutions with an efficient way to assess and assist students according to their academic performance. These models have the potential to revolutionize academic guidance and support, improving student outcomes in a data-driven setting in the process. The future of education is expected to be shaped by sophisticated machine learning techniques that prioritize accuracy, precision, recall, and F1-Score as the volume and complexity of educational data continue to rise.

## REFERENCES

[1] S. Hussain and M. Q. Khan, "Student-performulator: Predicting students' academic performance at secondary and intermediate level using machine learning," Annals of data science, vol. 10, no. 3, pp. 637–655, 2023.

[2] Vijayalakshmi and K. Venkatachalapathy, "Comparison of predicting student's performance using machine learning algorithms," International Journal of Intelligent Systems and Applications, vol. 11, no. 12, p. 34, 2019.

[3] M. Yağcı, "Educational data mining: prediction of students' academic performance using machine learning algorithms," Smart Learning Environments, vol. 9, no. 1, p. 11, 2022.

[4] J. Xu, K. H. Moon, and M. Van Der Schaar, "A machine learning approach for tracking and predicting student performance in degree programs," IEEE J Sel Top Signal Process, vol. 11, no. 5, pp. 742–753, 2017.

[5] Acharya and D. Sinha, "Early prediction of students performance using machine learning techniques," Int J Comput Appl, vol. 107, no. 1, pp. 37–43, 2014.

[6] D. M. Ahmed, A. M. Abdulazeez, D. Q. Zeebaree, and F. Y. H. Ahmed, "Predicting university's students performance based on machine learning techniques," in 2021 IEEE International Conference on Automatic Control & Intelligent Systems (I2CACIS), IEEE, 2021, pp. 276–281.

[7] H. Pallathadka, A. Wenda, E. Ramirez-Asís, M. Asís-López, J. Flores-Albornoz, and K. Phasinam, "Classification and prediction of student performance data using various machine learning algorithms," Mater Today Proc, vol. 80, pp. 3782–3785, 2023.

[8] H. Turabieh, "Enhanced Binary Genetic Algorithm as a Feature Selection to Predict Student Performance," 2021.

[9] Asselman, M. Khaldi, and S. Aammou, "Enhancing the prediction of student performance based on the machine learning XGBoost algorithm," Interactive Learning Environments, vol. 31, no. 6, pp. 3360–3379, 2023.

[10] H. Altabrawee, O. A. J. Ali, and S. Q. Ajmi, "Predicting students' performance using machine learning techniques," JOURNAL OF UNIVERSITY OF BABYLON for pure and applied sciences, vol. 27, no. 1, pp. 194–205, 2019.

[11] E. S. Bhutto, I. F. Siddiqui, Q. A. Arain, and M. Anwar, "Predicting students' academic performance through supervised machine learning," in 2020 International Conference on Information Science and Communication Technology (ICISCT), IEEE, 2020, pp. 1–6.

[12] M. Koutina and K. L. Kermanidis, "Predicting postgraduate students' performance using machine learning techniques," in International Conference on Engineering Applications of Neural Networks, Springer, 2011, pp. 159–168.

[13] Y. A. Alsariera, Y. Baashar, G. Alkawsi, A. Mustafa, A. A. Alkahtani, and N. Ali, "Assessment and evaluation of different machine learning algorithms for predicting student performance," Comput Intell Neurosci, vol. 2022, 2022.

[14] S. S. Shreem, H. Turabieh, S. Al Azwari, and F. Baothman, "Enhanced binary genetic algorithm as a feature selection to predict student performance," Soft comput, vol. 26, no. 4, pp. 1811–1823, 2022.

[15] P. Cortez and A. M. G. Silva, "Using data mining to predict secondary school student performance," 2008.

[16] J. L. Rastrollo-Guerrero, J. A. Gómez-Pulido, and A. Durán-Domínguez, "Analyzing and predicting students' performance by means of machine learning: A review," Applied sciences, vol. 10, no. 3, p. 1042, 2020.

[17] P. Dabhade, R. Agarwal, K. P. Alameen, A. T. Fathima, R. Sridharan, and G. Gopakumar, "Educational data mining for predicting students' academic performance using machine learning algorithms," Mater Today Proc, vol. 47, pp. 5260–5267, 2021.

[18] S. Hashim, W. A. Awadh, and A. K. Hamoud, "Student performance prediction model based on supervised machine learning algorithms," in IOP Conference Series: Materials Science and Engineering, IOP Publishing, 2020, p. 32019.

[19] Masci, G. Johnes, and T. Agasisti, "Student and school performance across countries: A machine learning approach," Eur J Oper Res, vol. 269, no. 3, pp. 1072–1085, 2018.

[20] Thammasiri, D. Delen, P. Meesad, and N. Kasap, "A critical assessment of imbalanced class distribution problem: The case of predicting freshmen student attrition," Expert Syst Appl, vol. 41, no. 2, pp. 321–330, 2014.

[21] K. Pal and S. Pal, "Data mining techniques in EDM for predicting the performance of students," International Journal of Computer and Information Technology, vol. 2, no. 06, pp. 2279–2764, 2013.

[22] Osmanbegovic and M. Suljic, "Data mining approach for predicting student performance," Economic Review: Journal of Economics and Business, vol. 10, no. 1, pp. 3–12, 2012.

[23] Kabakchieva, "Student performance prediction by using data mining classification algorithms," International journal of computer science and management research, vol. 1, no. 4, pp. 686–690, 2012.

[24] S. B. Kotsiantis, "Use of machine learning techniques for educational proposes: a decision support system for forecasting students' grades," Artif Intell Rev, vol. 37, pp. 331–344, 2012.

[25] M. Hussain, W. Zhu, W. Zhang, S. M. R. Abidi, and S. Ali, "Using machine learning to predict student difficulties from learning session data," Artif Intell Rev, vol. 52, pp. 381–407, 2019.

[26] Das, A. Stein, N. Kerle, and V. K. Dadhwal, "Landslide susceptibility mapping along road corridors in the Indian Himalayas using Bayesian logistic regression models," Geomorphology, vol. 179, pp. 116–125, 2012.

[27] M. Braik, M. H. Ryalat, and H. Al-Zoubi, "A novel meta-heuristic algorithm for solving numerical optimization problems: Ali Baba and the forty thieves," Neural Comput Appl, vol. 34, no. 1, pp. 409–455, 2022, doi: 10.1007/s00521-021-06392-x.

[28] P. Sharma, S. Thangavel, S. Raju, and B. R. Prusty, "Parameter Estimation of Solar PV Using Ali Baba and Forty Thieves Optimization Technique," Math Probl Eng, vol. 2022, p. 5013146, 2022, doi: 10.1155/2022/5013146.

[29] Kumar, V. Sharma, and R. Naresh, "Leader Harris Hawks algorithm based optimal controller for automatic generation control in PV-hydro-wind integrated power network," Electric Power Systems Research, vol. 214, p. 108924, 2023.

[30] M. K. Naik, R. Panda, A. Wunnava, B. Jena, and A. Abraham, "A leader Harris hawks optimization for 2-D Masi entropy-based multilevel image thresholding," Multimed Tools Appl, pp. 1–41, 2021.

# Design of Teaching Mode and Evaluation Method of Effect of Art Design Course from the Perspective of Big Data

Danjun ZHU[1], Gangtian LIU[2]*

School of ART and DESIGN, Henan University of Science and Technology, Henan Province 471000, China

*Abstract*—In modern educational curriculum teaching, we should fully leverage the advantages of modern technology, especially in teaching methods, and deeply understand and apply big data technology. This article explores the design and effectiveness evaluation methods of curriculum teaching models from the perspective of big data. We utilized big data thinking and conducted research and practical exploration to compare and evaluate teaching mode design methods. In the art and design course, we adopted a blended learning model, combining MOOC and SPOC, and innovated traditional teaching methods and plans. Meanwhile, we investigated the teaching effectiveness and feasibility of this blended learning model. By extensively evaluating teaching techniques, evaluation methods, and technologies that support the learning process, we reconstructed blended learning evaluation indicators and evaluated the effectiveness of learning outcomes and processes under different teaching modes. The research results show that the blended learning model based on big data perspective can significantly improve the effectiveness of classroom teaching. In contrast, learners' self-learning ability and practical innovation ability have also been further improved.

*Keywords—Big data perspective; teaching mode; evaluation system; art and design; hybrid teaching*

## I. INTRODUCTION

This throughout the data explosion of the 21st century, data is no longer a static and obsolete number [1] but has become a business capital. This important economic input can create new economic benefits. From its beginnings in business and technology, it is gradually moving towards and having a huge impact on healthcare and education. From the perspective of big data, the teaching mode design of art and design courses needs to fully consider the learning needs and practical ability cultivation of students. By collecting learning data from students and analyzing their learning behaviors and habits, teachers can develop teaching plans and plans that are more in line with their actual needs. For example, teachers can use data mining techniques to analyze students' learning trajectories, stay times, and review times, understand their learning difficulties and needs, and thus develop more accurate teaching plans.

In addition, in the process of designing teaching modes, it is also necessary to fully consider the setting of course structure and content. Art and design courses usually include two parts: theoretical knowledge and practical operation, and the application of big data technology can better promote the integration of these two parts. Through big data analysis, teachers can have a clearer understanding of students' learning situations and needs; thereby better adjusting course structure and content, and improving teaching effectiveness.

In such an era, the practice and exploration of art design course teaching [2] mode design and assessment and evaluation methods are an important part of which has a guiding role, an important way to strengthen the construction of courses, professional construction and improve the quality of teaching, as well as a method to test the quality of talent training in colleges and universities. In the era of big data, it is a question worth practicing and exploring how to draw on the big data thinking that has triggered social changes and carry out scientific reform for the assessment and evaluation methods of art and design courses around teaching purposes, student characteristics and course objectives of art and design majors in colleges and universities [3].

In order to break the status quo of subjectivity, one-sidedness, and singularity of assessment and evaluation methods, the design assessment and evaluation methods of art design course teaching mode in universities need to be reformed. The two cores of big data thinking are important to draw on in this process. In the "small data era", because the amount of information collected is relatively small, subtle errors will be magnified in the limited information, and the error rate will increase, so much so that it may affect the accuracy of the whole result. In assessing and evaluating the design of the teaching mode of art and design courses in colleges and universities, a single, one-sided quantitative evaluation standard cannot guarantee that the evaluation results are free from bias. Therefore, quantitative evaluation criteria in assessing and evaluating art and design courses can only be one component of the information and data required. The depth and breadth of data and information should be increased, and the course assessment and evaluation should be conducted with more comprehensive and diversified data and information to promote a more comprehensive assessment and evaluation [4].

Professional competencies such as method, thought, creativity, and expression in creating art and design works are diverse, as are the professional competencies involved, such as teamwork, re-learning, and resistance to frustration. These competencies should be emphasized and focused on in university art and design courses, and they can all be learned, exercised, and presented in the design process. Focusing on the learning and design process and conducting real-time

assessment and evaluation in stages for the corresponding competency points can prompt students to enrich their training in more professional and vocational competencies, clarify their learning and practice goals, and improve their initiative. In the meantime, the teacher is deeply involved in the whole process of students' learning and design. Based on what the students observed in the learning and creation process and the actual situation of learning competency points, the teacher can promptly adjust the teaching progress, supplement or strengthen the relevant contents, and give students adequate guidance.

As a practical and innovative discipline, art and design courses require the support of big data technology in order to better cultivate students' innovative thinking and practical abilities. From the perspective of big data, the teaching mode design of art and design courses needs to fully consider the learning needs and practical ability cultivation of students. By collecting learning data from students and analyzing their learning behaviors and habits, teachers can develop teaching plans and plans that are more in line with their actual needs. For example, teachers can use data mining techniques to analyze students' learning trajectories, stay times, and review times, understand their learning difficulties and needs, and thus develop more accurate teaching plans. Diversity of assessment and evaluation subjects to make assessment and evaluation more objective, the following points need to be achieved. Firstly, in course assessment and evaluation, it is usually necessary for the lecturer to take the lead and several non-lecturing professional teachers to participate. Secondly, there is a need to implement a three-in-one collaborative assessment in the assessment of courses. Thirdly, experts from academia and industry were invited to conduct immediate assessments of students to promote teaching and learning. Fourthly, elements of professional competitions in-course assessment and evaluation are introduced to get evaluated on more platforms. Teachers can guide students to choose professional competitions for university students with high gold content, wide audience, and indicative nature. They can also evaluate students and course works based on their performance and achievements in professional skills competitions.

Drawing on big data thinking, the method of diversifying assessment and evaluation subjects and assessment and evaluation content can effectively improve the problems of subjunctivization, one-sidedness, and singularity that arise in the assessment and evaluation of art and design courses in colleges and universities so that the assessment and evaluation of art and design courses is no longer a task at the end of the course, but a part of the course teaching in which students must participate. Using the method of diversifying assessment and evaluation subjects and assessment and evaluation contents not only makes the assessment and evaluation more comprehensive and objective, but more importantly, it makes the learning objectives clear to students while improving their learning initiative.

In recent years, MOOC has developed rapidly, with the construction of various MOOC platforms and a sudden surge in the number of online courses and users [5]–[7]. The huge growth in quantity has also caused certain quality problems, which are mainly manifested as follows: firstly, learners do not have a better grasp of their knowledge needs, leading to a large number of learners blindly registering and following the trend of learning, and MOOC platforms do not have effective means to restrain learners' learning discipline and learning progress. Secondly, teachers' lack of necessary guidance and communication during the learning process has reduced learners' interest in learning. Although most MOOC platforms have some interactive means, they cannot find learners' questions and answer them promptly. The one-way knowledge dissemination method based on video learning is unsuitable for cultivating and improving learners' learning capabilities. It is even more difficult to achieve in-depth learning. Thirdly, the MOOC platform only provides some theoretical materials and practice content without monitoring the learning process of learners. A new evaluation system should be developed by proposing a new teaching model for art and design courses from the perspective of big data.

The concept of SPOC [8]–[10] (Small Private Online Course) was first introduced and used by Professor Armand Fox of California State University, Berkeley. Massive and Open are in opposition to small and private. While Small helps to increase participation, interactivity, and completion rates, Private make the course somewhat restrictive and simple to keep up and manage. The hybrid teaching mode, which combines the ideas of MOOC and the flexibility of SPOC, is a further learner-centred approach that takes into account the leading role of the teacher, incorporates multiple teaching methods, mixes multiple teaching devices, combines multiple learning resources, and uses multiple evaluation indicators to build a hybrid learning mode that intersects synchronous and asynchronous learning [11], [12].

The construction of a teaching mode for art and design courses in higher education is upgraded through the borrowing and application of the big data thinking mode while forming a diversified assessment work system, both in terms of assessment content and subject matter respectively. The result is that while eliminating the problem of developmental limitations, the overall art course assessment can be better adapted to student growth needs and encourages and actively contributes to the advancement of art and design education. This paper aims to analyze the effectiveness of teaching methods and the feasibility of this mode by using a hybrid teaching mode that combines MOOC and SPOC. By conducting a comprehensive evaluation of teaching techniques, assessment methods and technologies that support the learning process, and a reconstructed hybrid teaching assessment index to assess learning outcomes in different teaching modes, the design of a teaching mode for art and design courses from a big data perspective act as a catalyst for traditional teaching methods.

In previous studies, scholars mainly focused on the potential applications and advantages of big data in art and design courses. For example, some studies suggest that big data can help teachers better understand students' learning needs and behavioral patterns, thereby developing more personalized teaching plans. In addition, research has shown that big data can be effectively used to evaluate the learning effectiveness of students and the effectiveness of course design.

In the current research, we will delve deeper into the design of teaching modes and evaluation methods for art and design courses from the perspective of big data. Firstly, we will further investigate and address the challenges of data processing and analysis, and propose more effective methods and techniques to handle large amounts of unstructured data. Secondly, we will delve deeper into privacy and security issues and propose more comprehensive data protection strategies and measures. In addition, we will also explore how to combine big data with other advanced technologies such as artificial intelligence; machine learning, etc. to provide a more personalized and efficient learning experience.

## II. MAIN FRAMEWORK FOR BLENDED LEARNING

The SPOC-based hybrid teaching combines traditional face-to-face and independent learning on the SPOC platform. Considering the actual teaching mode design of art design courses in colleges and universities, this paper divides the whole learning stage into task-based independent learning before class, guided learning through teacher-student interaction within the class, and enhanced learning through consolidation and evaluation after class. The teaching framework of the blended teaching mode is shown in Fig. 1.

The teacher sets the learning tasks before the class; learners learn through online videos, share their learning, and practice through online tests to consolidate and master the knowledge. Teachers led workshops in the classroom, practiced and analyzed the lessons, refined knowledge, and reviewed learning outcomes. At the end of the lesson, the teacher assigns practical tasks, and learners expand their knowledge, deepen their understanding and mastery, and submit their results to the teacher for comments [13], [14].

Learners can access theoretical and practical knowledge through various channels, depending on their preferences. They can use materials such as textbooks, online videos, electronic lesson plans, and supplementary materials provided in the online courses. In addition, learners can use the case library to practice their art and design expertise. Its components are shown in Fig. 2.

*1)* Includes 56 videos on course knowledge and 18 videos on extended knowledge for learners' independent study before class and extended practice after class.

*2)* 38 videos of experimental explanations and 42 experimental assignments, and learners can use the online assessment platform to program and receive real-time feedback.

*3)* Five unit-based theory tests, five unit-based arts, design tests and a final theory test and a practical test.

*4)* From the fifth week of the course, learners will work in groups to solve two complete cases. Rigorous tests and an assessment of the system's hand-in results are conducted as required.



Fig. 1. MOOC + SPOC-based hybrid teaching mode with a modal teaching framework.



Fig. 2. Art and design course hybrid teaching plan.

## III. Hybrid Teaching Evaluation Indicators

In the hybrid teaching approach, the assessment of learners focuses heavily on the entire learning process. It involves monitoring their pre-learning progress, assessing their mastery of knowledge during class, and measuring their ability improvement after class. This solves the problem of traditional teaching methods neglecting to assess learners' learning attitudes and learning processes. Through the evaluation system, teachers can identify the shortcomings of teaching strategies and content, understand learners' learning, improve teaching design in time, and make timely and targeted interventions on learners' learning to improve teaching quality.

In this study, the Delphi method [15], [16], the Experts Grading Method (EGM), and the Analytic Hierarchy Process (AHP) [17], [18] were combined to construct the index system, and 16 industry experts were sent a call for comments to establish a hierarchical mode of the evaluation system, as shown in Fig. 3. The weights for each indicator system were determined and presented in Table I.

Fig. 4 compares the total weights of the blended learning assessment indicators. As can be seen from Fig. 4, the indicator system no longer uses examination results as the only criterion for evaluating learners' abilities. Still, it gives more weight to problem-solving, active learning, participation in discussions,

and design innovation, thus placing more emphasis on cultivating independent learning and practical and innovative abilities in the blended teaching mode and more emphasis on exploring and expanding learners' potential abilities.

### A. Delphi Method

*1) Traditional delphi method:* The Delphi method emerged in the early 1950s as a predictive technique invented by Dalktey and his associates and has been widely used in curriculum teaching. By conducting a questionnaire for a decision-making group, not only can a brainstorming effect be achieved, but it can also be revised repeatedly to obtain a final result. The characteristics of the Delphi method are as follows.

*a)* The Delphi method relies on the experience and judgment of the participants, and the intervention of individual subjectivity is inevitable but is, therefore, fully inclusive of a diversity of views.

*b)* The different participants participate in the analysis of the thesis anonymously to avoid human interference.

*c)* The indicators of prediction and judgment in the questionnaire need to be studied carefully and focus on the feedback of participants' opinions. The final results will converge to reach a consensus through the analysis of the questionnaire results and repeated surveys.



Fig. 3. Evaluation index system of the blended teaching mode.

TABLE I. WEIGHTING OF EVALUATION INDICATORS FOR BLENDED LEARNING AND TEACHING

| First-level indicator | Weights | Secondary indicators | Weights | Combined weights |
|---|---|---|---|---|
| learning attitude | 0.128 | Number of logins to the platform | 0.151 | 0.019 |
| | | watch video time | 0.524 | 0.067 |
| | | Posted views | 0.325 | 0.042 |
| learning ability | 0.234 | preview course | 0.456 | 0.107 |
| | | Questions in class | 0.544 | 0.127 |
| Practical ability | 0.234 | Extracurricular Reading | 0.264 | 0.062 |
| | | problem-solving skills | 0.736 | 0.250 |
| Academic record | 0.404 | Chapter test | 0.413 | 0.167 |
| | | Final assessment | 0.587 | 0.237 |

Fig. 4.   Comparison of composite indicator weights for blended learning.

*2) Fuzzy Delphi Method (FDM):* As the Delphi method is used to make decisions through multiple questionnaires and integrate expert opinions, its conclusions are more rigorous and reasonable, so it is more widely used. However, it has drawbacks, such as poor questionnaire design and convergence due to widely differing expert opinions. The Delphi method may ignore the correct and unique opinions of a few experts, and the method expresses expert opinion in precise numerical terms, which is not in line with the ambiguity of human thinking. The steps are as follows:

*a)* Building set of influencing factors.

*b)* Summarizing expert opinion

Based on the different influencing factors obtained from the aggregation, a questionnaire was designed, and an expert questionnaire was administered. Considering human thinking judgments

To make the expert's opinion complete, the expert is asked to fill in the questionnaire, according to his professional knowledge, to determine the degree of influence of each influencing factor on the selection and to give a value between 1 (no influence) and 10 (absolute influence), the higher the score means the greater the influence, the maximum value of the interval represents the maximum influence of the factor in the expert's opinion. In contrast, the interval's minimum value represents the factor's minimum influence on the expert's opinion. The range's maximum value represents the factor's maximum influence as perceived by the expert. In contrast, the range's minimum value represents the factor's minimum influence as perceived by the expert.

*3)* Establishing trigonometric fuzzy functions [19] and integrating expert opinion.

The expert opinions obtained from the questionnaire were collated and organized.

$$\underset{\sim}{n} = (l_A, m_A, u_A) \tag{1}$$

$$\underset{\sim}{N} = (L_A, M_A, \quad U_A) \tag{2}$$

Create fuzzy sets of the minimum and maximum influence of each influence factor and, $\underset{\sim}{N}$ respectively, where

$l_A: min(x_{Ai}), \quad i = 1 \sim n$

$m_A: (x_{A1}, x_{A2}, \dots \dots x_{An})1/n, \quad i = 1 \sim n$

$u_A: max(xAi), \quad i = 1 \sim n$

$L_A: min(X_{Ai}), i = 1 \sim n$

$M_A: (X_{A1}, X_{A2}, \dots \dots X_{An})^{1/n}, i = 1 \sim n$

$U_A: max(X_{Ai}), \quad i = 1 \sim n$

$X_{Ai}$: The value of *ith* expert's maximum influence on the *A* influencing factor.

$L_A$: Lower limit for the value of the maximum influence of the expert group on the impact factor *A*.

$M_A$: Geometric mean of the expert group's assessment of the maximum impact of influence *A*.

$U_A$: The upper limit of the value of the maximum influence of the expert group on the impact factor *A*.

Then, the affiliation functions of the fuzzy sets $\underset{\sim}{n}$ and $\underset{\sim}{N}$ can be expressed respectively as,

$$\mu_n(x_i) = \begin{cases} 0 & 0 < x < l_A \\ \dfrac{x - l_A}{m_A - l_A} & l_A \leq x < m_A \\ \dfrac{-x + u_A}{u_A - m_A} & x = m_A \\ x > m_A \end{cases} \tag{3}$$

$$\mu_N(x_i) = \begin{cases} 0 & 0 < x < L_A \\ \dfrac{x - L_A}{M_A - L_A} & L_A \le x < M_A \\ 1 & x = M_A \\ \dfrac{-x + U_A}{U_A - M_A} & x > M_A \end{cases} \qquad (4)$$

*4) Calculate the quantitative value of each influencing factor:* If the opinions of the interviewed experts have converged, i.e. $u_A \ge L_A$ , and $G_A < Z_A$ , then the intersection of the two fuzzy sets $CP$ can be found as the quantitative value of this influence factor, which represents the degree of influence of this influence factor on the target as agreed by the experts as a whole, and can be used as the basis for screening the influence factor. From the above figure, the intersection point $CP$ can be solved by the following equation.

$$\frac{-x + u_A}{u_A - m_A} = \frac{x - L_A}{M_A - L_A} \qquad (5)$$

$$CP = x = \frac{u_A * M_A - m_A * L_A}{M_A - L_A + u_A - m_A} \qquad (6)$$

*5) Set domain values and filter influencing factors:* According to the principle of setting indicators for equipment selection, set suitable indicators and indicator filtering domain values $\alpha$, and filter the $CP$ values of each influencing factor obtained above, with the rule that if $CP \ge \alpha$, the influencing factor will be used as a criterion or indicator, and if $CP < \alpha$, the influencing factor will be ignored.

The advantage of using the fuzzy Delphi method is that integrating expert opinion with fuzzy functions better expresses the vagueness of human thinking and the lack of Certainty. The overlap region is checked to determine whether the expert opinion has converged, making the analysis more rigorous and reasonable. Using the concept of 'maximum and minimum values within the range of possibilities' instead of 'most probable' and 'least probable' enhances the usefulness of the Delphi method.

*B. Hierarchical Analysis*

Hierarchical analysis is a method that combines quantitative and qualitative approaches by organizing and synthesizing the opinions of people's subjective judgments. It is also effective in dealing with complex problems difficult to analyze by quantitative methods, breaking down complex problems into levels for step-by-step analysis [20]. The method allows one's subjective judgment to be expressed and processed in quantitative form and can also suggest whether one's subjective judgment on a particular type of problem is inconsistent. The method is now widely used on multi-objective optimization problems and can determine the weights of individual objectives, thus assisting in decision-making. With the rapid development of technology, deep learning, as an important branch of artificial intelligence, has gradually penetrated into every corner of schools [21]. The introduction of this technology not only brings new teaching methods to education, but also poses challenges to traditional educational models. How to understand and grasp the relationship between deep learning and education has become a topic that we need to explore in depth at present. The entry of deep learning into schools signifies a change in teaching methods. The traditional education model often centers on teachers, while deep learning emphasizes student-centered learning, utilizing artificial intelligence technology to provide personalized learning experiences [22]. Through the analysis of a large amount of data, deep learning can accurately grasp the learning needs and habits of students, and provide teachers with more scientific teaching suggestions. At the same time, it can also help students better understand knowledge and improve learning efficiency. However, the introduction of deep learning has also brought some challenges. Firstly, data privacy and security issues cannot be ignored. During the process of using deep learning technology, students will generate a large amount of data, and how to ensure the privacy and security of this data has become a major challenge. Secondly, the effective application of deep learning technology requires teachers to possess corresponding technical literacy, which is a significant challenge for many teachers [23]. In addition, excessive reliance on technology may lead to the neglect of humanistic care in education, which is a problem that we need to be vigilant about while pursuing technological progress.

In brief, the hierarchical approach begins with a description of the problem, followed by identifying the influencing factors and establishing a hierarchy of relationships. The relative importance of the decision factors at each level is identified by using pairwise comparisons on a scale of proportionality, from which positive and negative comparisons are established, calculating the eigenvalues and eigenvectors of the matrix, and finding the weights of each attribute, the important steps are explained below.

*1) Description of the problem:* When conducting an AHP, the system in which the problem is situated should be analyzed in as much detail as possible, with all the factors that may affect the problem being included in the problem and the main objectives of the problem being determined, but with attention to the interrelationship and independence of the factors.

*2) Creating hierarchy:* The interaction of many factors influences the selection of equipment. This study uses a logical thinking approach to consider the factors that may affect equipment selection, divides the criteria and indicators for evaluation into different levels of varying importance, and then examines the following levels, starting with the highest level of objectives.

*3) Building a judgment matrix:* The judgment matrix is created by using one of the factors at a higher level in the hierarchical mode as an evaluation criterion, and the experts make a two-by-two comparison of the factors at this level, using a judgment scale to determine the matrix elements. To assess the relative importance of indicators, a common practice is to employ the nine-point scale ranging from 1 to 9. If a criterion has $n$ factors at the lower level, the judgment matrix $F_1$、$F_2 \cdots F_n$ $A$ is created, as shown in Fig. 7.

$$A = \begin{bmatrix} 1 & a_{12} & \cdots & a_{1n-1} & a_{1n} \\ 1/a_{12} & 1 & \cdots & a_{2n-1} & a_{2n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1/a_{1n-1} & 1/a_{2n-1} & \cdots & 1 & a_{n-1n} \\ 1/a_{1n} & 1/a_{2n} & \cdots & 1/a_{n-1n} & 1 \end{bmatrix} \tag{7}$$

Making use of the judgment matrix $A$ to determine each factor's weight at each level concerning the corresponding elements from the level before, i.e., calculate the highest attribute root $\lambda_{max}$ of $A$, matching the normalized value eigenvector $v$, i.e.

$$Av = \lambda_{max} \tag{8}$$

and

$\sum_{i=1}^{n} v_i = 1$, where $v = (v_1, v_2, \cdots, v_n)^T$.

To achieve a scientific and objective calculation of each factor's weight, it becomes essential to establish the equilibrium of the judgment matrix A. The judgment matrix's consistency assessment is measured by the indicator known as,

$$CI = \frac{\lambda_{max}}{n-1} \tag{9}$$

As the size of the index of consistency ($CI$) increases, the consistency of the judgment matrix worsens. On the other hand, when $CI$ equals zero, it indicates complete satisfaction with the consistency of the judgment matrix.

Consistency test coefficient

$$CR = \frac{CI}{RI} \tag{10}$$

where, $RI$ is the random consistency indicator associated with the order of the judgment matrix, and its correspondence can be found in Table II.

The correlation between the judgment matrix's order and the random consistency index is shown in Fig. 5, and it shows that an increase in the corresponding random consistency index accompanies an increase in the judgment matrix's order.

When $CR < 0.1$, the consistency test is passed by the judgment matrix, when $CR \geq 0.1$, the judgment matrix fails to meet the consistency test and thus requires correction.



Fig. 5. Plot of the order of the judgment matrix against the random consistency index.

TABLE II. RANDOM CONSISTENCY INDICATORS

| n | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| RI | 0.05 | 0.23 | 0.46 | 0.80 | 1.22 | 1.29 | 1.35 | 1.48 | 1.54 |

## IV. ANALYSIS OF THE EFFECTIVENESS OF THE BLENDED TEACHING MODE

### A. Analysis of the Learning Process

This article uses online platforms to publish questionnaires and extensively collect data from different regions and types of schools. Go deep into specific universities, have face-to-face communication with teachers and students, and obtain first-hand information. Utilize the school's academic management system and other database resources to obtain data on course teaching modes, student grades, and other aspects. To provide a reasonable analysis of the blended teaching model, 226 students were selected from the art and design majors of the China Academy of Art. By dividing these students into two groups, both groups' basic theoretical knowledge and practical design ability were tested, and the experimenter's entry grades were determined based on the test results. The first group consisted of 112 students with an average entry score of 159.56, including an average score of 103.84 in theory and 183.63 in practical design; the second group consisted of 114 students with an average entry score of 158.73, including an average score of 103.21 in theory and 183.33 in practical design, which was generally consistent between the two groups. Group 1 used a hybrid teaching method, and Group 2 used a traditional lecture method. Both groups' teaching hours were 48 hours, as shown in Table III.

A visual comparison of the entry scores of the two groups of students is represented in Fig. 6. It is clear from Fig. 6 that the difference between the entry scores of the two groups of students is not significant. However, the scores of all indicators of the blended teaching mode are higher than those of the traditional teaching, which shows the superiority of the blended teaching mode.

An independent samples t-test was first conducted for each of the two groups of students' scores, and the results are shown in Table IV.

TABLE III. COMPARISON OF ENTRY SCORES OF STUDENTS IN THE TWO GROUPS

| Category | Entrance grades | Theoretical score | Design in Practice |
|---|---|---|---|
| First group | 159.56 | 103.84 | 183.64 |
| Second Group | 158.73 | 103.21 | 183.33 |

Fig. 6. Visual comparison of the entry scores of the two groups of students.

TABLE IV.    T-TEST FOR INDEPENDENT SAMPLES OF STUDENTS' PERFORMANCE IN EACH OF THE TWO TEACHING STYLES

| Category | blended teaching | traditional teaching | T | Sig |
|---|---|---|---|---|
| Overall review | 82.77 | 78.14 | 2.224 | 0.030 |
| practice | 47.02 | 42.33 | 2.868 | 0.150 |
| theory | 31.84 | 32.64 | 1.199 | 0.234 |

A visual comparison of the three-sample t-tests of achievement in the blended and traditional teaching modes in the art and design course is shown in Fig. 7.



(a) Teaching mode comparison.



(b) T and Sig

Fig. 7. Independent sample T-test for two student grades (a) Teaching mode comparison (b) T and Sig.

*1)* Regarding the overall assessment results, the standard deviation of students' performance under the traditional way of teaching was 11.594, indicating a wide range of performance distribution indicating the existence of polarization of students. Meanwhile, the t-test result is: $t = 2.224, P = 0.03 < 0.05$, indicating that there is a significant difference between the two teaching effects. Hence, it can be inferred that the hybrid teaching mode is considerably more effective in terms of teaching than the traditional teaching method.

*2)* Regarding practical performance, the independent sample t-test result indicates a significant difference in practical operation between students of hybrid teaching and traditional teaching mode. It indicates that the hybrid teaching mode is more helpful to the exercise and improvement of students' practical skills.

*3)* The results of the independent samples t-test from the test paper quiz scores were: $t = 1.119, P = 0.0.234 > 0.05$, indicating no significant difference in theoretical knowledge learning between the students in the two modes of teaching.

### B. Evaluation of the Learning Process

To analyze the effect of the blended teaching approach and understand the students' experience and feelings towards the teaching organization. A questionnaire was designed and administered to the experimenters, which contained questions on interest in learning art and design courses, knowledge acquisition, learning attitudes and motivation, knowledge transfer and application, and improvement of abilities, each

containing multiple questions. The analysis and summary of the survey results are as follows:

*1)* The data on students' interest in studying art and design courses before and after the course was offered are shown in Tables 5 and 6, respectively.

A comparison of the data in Tables V and VI is shown in Fig. 8. It can be seen from Fig. 8 that student interest is significantly increased under blended learning.

*2)* The survey counted the seriousness of students' learning under the two modes and investigated the degree of influence of learning attitudes on practical ability. The results indicated that when students' seriousness of learning was

comparable under the two modes. Still, there was a certain difference in practical ability; students with the blended teaching model had a certain advantage in practical ability, indicating that blended teaching was more helpful in improving students' practical ability. The result of the survey shows that most students are more open, diverse, and innovative in their problem-solving when using the blended teaching mode.

*3)* The ability of students who adopted the blended teaching mode is also investigated regarding knowledge transfer. The specific data are shown in Table VII.

TABLE V.    DATA ON STUDENTS' INTEREST IN STUDYING ART AND DESIGN COURSES BEFORE THE COURSE WAS OFFERED

| Before class starts | very interested | interested | not interested |
|---|---|---|---|
| the first time | 44 | 120 | 12 |
| the second time | 55 | 110 | 11 |
| the third time | 40 | 118 | 18 |
| the fourth time | 50 | 115 | 11 |
| the fifth time | 45 | 112 | 19 |

TABLE VI.    DATA ON STUDENTS' INTEREST IN STUDYING ART AND DESIGN COURSES AFTER THE COURSE WAS OFFERED

| End of the school term | very interested | interested | not interested |
|---|---|---|---|
| the first time | 84 | 90 | 2 |
| the second time | 80 | 93 | 3 |
| the third time | 76 | 96 | 4 |
| the fourth time | 80 | 92 | 4 |
| the fifth time | 78 | 93 | 5 |



(a) Before class starts.

(b) End of the school term.

Fig. 8.   Data on students' interest in learning art and design courses before and after the course was offered (a) Before class starts (b) End of the school term.

TABLE VII.    DATA ON KNOWLEDGE TRANSFERABILITY BETWEEN THE BLENDED AND TRADITIONAL MODES OF TEACHING AND LEARNING

| Category | Ability to transfer knowledge | No knowledge transfer | Practical tasks can be completed | Practical tasks not completed |
|---|---|---|---|---|
| blended teaching mode | 61 | 51 | 69 | 43 |
| traditional teaching mode | 42 | 72 | 47 | 67 |

Fig. 9. Pie chart of knowledge transferability and task completion (a) Blended teaching mode (b) Traditional teaching mode.

A comparison of the knowledge transfer ability data and task completion pie charts for the two teaching modes is shown in Fig. 9. The results show that 54.1% of students in the blended mode could make connections between old and new knowledge, and 61.3% of students reported that they could complete the practical tasks as required. This figure is much higher than that of the students using the traditional teaching method (36.5% and 41.2% respectively). The difference between the two modes of teaching is not entirely a difference in the amount of knowledge acquired but rather a difference in the ability to link theory and practice and to connect old and new knowledge. This also shows that the mixed teaching mode is more conducive to transferring knowledge and cultivating practical problem-solving skills.

*4)* In terms of learning motivation, the survey found that 85.3% of students would use their spare time for pre-study and revision when using the blended teaching mode and 67.4% would actively consult extra-curricular materials related to the course, which is much higher than those using the traditional teaching mode (64.6% and 43.8% respectively). This indicates that the willingness and ability of students to learn have been significantly improved by using the blended teaching mode.

In addition, the responses to the open-ended questions in the questionnaire showed that students who adopted the blended teaching mode identified more strongly with the new teaching mode and generally felt that their interest in learning, confidence, and learning ability had increased significantly.

The analysis of students' performance and questionnaires reveals that the blended teaching mode is more likely to stimulate students' interest in learning, enhance their awareness of independent learning and learning ability, promote the acquisition, transfer, and application of knowledge, and improve their practical skills and problem-solving abilities.

In terms of effectiveness evaluation, the application of big data principles makes the evaluation of teaching effectiveness more objective, accurate, and comprehensive. By comparing the learning outcomes and processes under different teaching modes, teachers can have a clearer understanding of the advantages and disadvantages of the new teaching mode. At the same time, the application of correlation analysis and predictive analysis also makes the evaluation of teaching effectiveness more scientific, refined, and forward-looking, providing teachers with more personalized and targeted teaching suggestions. In addition, we also found that integrating the principles of big data into teaching mode design and effectiveness evaluation requires teachers to have corresponding technical literacy and abilities. Therefore, strengthening the training and guidance of teachers is one of the key to achieving the design of teaching modes and evaluation methods for art and design courses from the perspective of big data.

In summary, the design of teaching modes and evaluation methods for art and design courses from the perspective of big data has important practical significance and application value. By deeply exploring and analyzing the application of big data principles in teaching mode design and effectiveness evaluation, we can better understand students' learning needs and habits, develop more scientific, accurate, and personalized teaching plans, and improve teaching effectiveness and learning experience. Meanwhile, strengthening the technical literacy and ability cultivation of teachers is also an important guarantee for achieving this goal.

## V. DISCUSSION

From the perspective of big data, the teaching mode design of art and design courses needs to fully consider the learning needs and practical ability cultivation of students. By collecting learning data from students and analyzing their learning behaviors and habits, teachers can develop teaching plans and plans that are more in line with their actual needs. For example, teachers can use data mining techniques to analyze students' learning trajectories, stay times, and review times, understand their learning difficulties and needs, and thus develop more accurate teaching plans. The blended learning mode requires corresponding technical support, including network platforms, teaching software, etc. However, some schools are unable to provide sufficient technical support due to funding, technology, and other reasons, leading to limitations in the implementation of blended learning models. The blended learning model requires teachers to have corresponding teaching design and organizational abilities, as well as a certain level of information technology skills. However, some teachers find it difficult to

effectively implement blended learning models due to a lack of relevant experience and skills. The blended learning model requires active participation and cooperation from students. However, some students lack learning motivation and self-discipline, resulting in low participation and poor teaching effectiveness.

Schools should strengthen the training and guidance of teachers, improve their teaching design and organizational abilities, and cultivate their information technology skills in order to better implement blended learning models. Teachers should take effective measures to increase student engagement, such as setting interesting learning tasks, providing personalized learning support, etc., to stimulate students' interest and motivation in learning.

## VI. CONCLUSION

This paper explores the design of teaching modes and methods for art design courses using massive data. With massive data thinking, a comparative evaluation study assesses different curriculum teaching methods. A blended teaching mode combining MOOC and SPOC is adopted for the art design course to innovate the traditional teaching mode and plan. The study's results indicate that the utilization of blended teaching mode can potentially enhance student's performance and knowledge transfer ability and foster an improved learning motivation. Simultaneously, the awareness and learning capability of students' autonomous learning has also been enhanced, which promotes the acquisition, transfer, and application of knowledge and enhances students' practical skills and problem-solving capabilities. The teaching method and evaluation system proposed in the current work also have many areas that need to be revised and improved. However, the practice has proved that the blended teaching mode has great advantages in improving the effect of classroom teaching, improving the learners' autonomous learning ability, practical innovation ability, and sustainable development ability. It is believed that the blended teaching mode will become one of the important learning methods in the future.

Although big data provides a wealth of information, how to effectively and accurately process and analyze this data remains a challenge. Especially when dealing with unstructured data such as text comments, student works, etc., data cleaning, annotation, and mining require a lot of time and manpower. Privacy and security issues cannot be ignored in the process of collecting, storing, and using student data. How to ensure the security, compliance, and anonymity of data, prevent data leakage and abuse, is a challenge that big data must face when applied in the field of education.

In order to overcome the limitations of insufficient infrastructure, future research can focus on developing more efficient and low-cost big data storage and analysis technologies, providing better technical support for the education field. With the increasing prominence of data security and privacy issues, future research should focus more on how to effectively utilize data for educational analysis and evaluation while ensuring data security and compliance.

## REFERENCES

[1] B. M. Kraemer, "Rethinking discretization to advance limnology amid the ongoing information explosion," *Water Res*, vol. 178, p. 115801, 2020. https://doi.org/10.1016/j.watres.2020.115801.

[2] Q. Xie, "Practice, Immersion and Collaboration—the Teaching Exploration of the Integration of Traditional Culture into Art Design Course from the Constructivism Theory," in *the 6th International Conference on Arts, Design and Contemporary Education (ICADCE 2020)*, Atlantis Press, 2021, pp. 301–304. https://doi.org/10.2991/assehr.k.210106.058.

[3] J. Huang, "Path analysis of university education management based on big data technology," in *Journal of Physics: Conference Series*, IOP Publishing, 2021, p. 042087. https://doi.org/10.1016/j.im.2018.12.003.

[4] Alt, D., & Raichel, N. Problem-based learning, self-and peer assessment in higher education: towards advancing lifelong learning skills. *Research Papers in Education*, 37(3), 370-394, 2022. https://doi.org/10.1080/02671522.2020.1849371.

[5] Zhu, B., Zheng, Y., Ding, M., Dai, J., Liu, G., & Miao, L. (2023). A pedagogical approach optimization toward sustainable architectural technology education applied by massive open online courses. Archnet-IJAR: International Journal of Architectural Research, 17(3), 589-607. https://doi.org/10.1108/ARCH-07-2022-0151.

[6] Q. Tan, X. Yan, Z. Qin, and B. Fang, "Study on MOOC's English Learning Mode Based on Pattern Recognition," *Wirel Commun Mob Comput*, vol. 2022, 2022. https://doi.org/10.1155/2022/4814658.

[7] Y. Fan, J. Jovanović, J. Saint, Y. Jiang, Q. Wang, and D. Gašević, "Revealing the regulation of learning strategies of MOOC retakers: A learning analytic study," *Comput Educ*, vol. 178, p. 104404, 2022. https://doi.org/10.1016/j.compedu.2021.104404.

[8] A. E. Garzón, T. S. Martínez, M. L. J. Ortega, J. A. Marin, G. G. Gomez. Teacher training in lifelong learning—The importance of digital competence in the encouragement of teaching innovation. Sustainability, 12(7), 2852, 2020. https://doi.org/10.3390/su12072852.

[9] Y. Liao, Q. Huang, C. Wang, Z. Zuo, Y. Wang, and Q. Yu, "Knowledge graph and its applications in MOOC and SPOC," in *2019 2nd International Conference on Contemporary Education and Economic Development (CEED)*, 2019, pp. 301–305. https://doi.org/10.26914/c.cnkihy.2019.037305.

[10] Zhang, J., Sziegat, H., Perris, K., & Zhou, C. (2019). More than access: MOOCs and changes in Chinese higher education. Learning, Media and Technology, 44(2), 108-123. https://doi.org/10.1080/17439884.2019.1602541.

[11] Li, X., & Ji, Y. (2023). Creative Teaching: The Realization of Ideological and Political Education Based on Chinese Traditional Culture. Chinese Studies, 12(4), 373-389. https://doi.org/10.4236/chnstd.2023.124027.

[12] X. Yang, "Exploration and practice of online and offline mixed teaching in functional experiment teaching," *International Journal of Social Science and Education Research*, vol. 4, no. 10, pp. 27–31, 2021. https://doi.org/10.1109/ICISE-IE53922.2021.00030.

[13] L. Jiang, "Factors influencing EFL teachers' implementation of SPOC-based blended learning in higher vocational colleges in China: A study based on grounded theory," *Interactive Learning Environments*, pp. 1–20, 2022. https://doi.org/10.1080/10494820.2022.2100428.

[14] Li, Z., & Jiang, W. (2022). Research on the Teaching Reform of Inorganic Chemistry Based on SPOC and FCM during COVID-19. Sustainability, 14(9), 5707. https://doi.org/10.3390/su14095707.

[15] L. Shen, J. Yang, X. Jin, L. Hou, S. Shang, and Y. Zhang, "Based on Delphi method and analytic hierarchy process to construct the evaluation index system of nursing simulation teaching quality," *Nurse Educ Today*, vol. 79, pp. 67–73, 2019. https://doi.org/10.1016/j.nedt.2018.09.021.

[16] Zhang, F., Liu, Q., & Zhou, X. (2022). Vitality evaluation of public spaces in historical and cultural blocks based on multi-source data, a case study of Suzhou Changmen. Sustainability, 14(21), 14040. https://doi.org/10.3390/su142114040.

[17] W.-C. Yang, J.-B. Ri, J.-Y. Yang, and J.-S. Kim, "Materials selection criteria weighting method using analytic hierarchy process (AHP) with simplest questionnaire and modifying method of inconsistent pairwise comparison matrix," *Proceedings of the Institution of Mechanical Engineers, Part L: Journal of Materials: Design and Applications*, vol. 236, no. 1, pp. 69–85, 2022. https://doi.org/10.1177/14644207211039912.

[18] Y. Y. Cho and H. Woo, "Factors in evaluating online learning in higher education in the era of a new normal derived from an Analytic Hierarchy Process (AHP) based survey in South Korea," *Sustainability*, vol. 14, no. 5, p. 3066, 2022. https://doi.org/10.3390/su14053066.

[19] G. Qiu, "Application of Wavelet Packet and Fuzzy Algorithm in Power System Short Circuit Fault Classification," *Math Probl Eng*, vol. 2022, 2022. https://doi.org/10.1155/2022/2882456.

[20] X. Chen, Y. Fang, J. Chai, and Z. Xu, "Does intuitionistic fuzzy analytic hierarchy process work better than analytic hierarchy process?" *International Journal of Fuzzy Systems*, vol. 24, no. 2, pp. 909–924, 2022. https://doi.org/10.1007/s40815-021-01163-1.

[21] W. C. Gresse, H. C. Hauck, F. S. Pacheco, B. M. Bertonceli Bueno. Visual tools for teaching machine learning in K-12: A ten-year systematic mapping. Education and Information Technologies, 26(5), 5733-5778, 2021. https://doi.org/10.1007/s10639-021-10570-8.

[22] Wahab, O. A., Mourad, A., Otrok, H., & Taleb, T. (2021). Federated machine learning: Survey, multi-level classification, desirable criteria and future directions in communication and networking systems. IEEE Communications Surveys & Tutorials, 23(2), 1342-1397. https://doi.org/10.1109/COMST.2021.3058573.

[23] Perrotta, C., & Selwyn, N. (2020). Deep learning goes to school: Toward a relational understanding of AI in education. Learning, Media and Technology, 45(3), 251-269. https://doi.org/10.1080/17439884.2020.1686017.

# Research on Evaluation and Improvement of Government Short Video Communication Effect Based on Big Data Statistics

Man Xu

Department of Journalism, Communication University of China, Qiqihar University
Beijing 100020, China, Heilongjiang 161000, China

*Abstract*—**Mainstream media is no longer the only way for people to obtain information, and the official media no longer has absolute control. People can choose the form and content of receiving information according to their preferences, which poses a new challenge to the government departments that have always been serious. From the beginning of short video to its prosperity, the government has shown great interest in its characteristics and functions. It has started to layout short video of government affairs on platforms such as Tiktok and Kwai, opened accounts one after another, and actively participated in the production and dissemination of content. Through the continuous launch of well-designed "hot money", the popularity of government affairs short videos on Tiktok and other platforms continued to rise, harvested a large number of fans, attracted social attention, and also brought good results and repercussions. This paper proposes an optimization design scheme for the evaluation and improvement of the dissemination effect of government short video based on big data statistics. The basic situation of government video is obtained through content analysis, and then the judgment coefficient and linear regression in big data statistics are used to extract common factors to improve the dissemination effect of government short video, so as to improve the dissemination influence of government short video. Finally, simulation test and analysis are carried out. Simulation results show that the proposed algorithm has certain accuracy, which is 8.24% higher than the traditional algorithm. Carrying out the research on the promotion planning and design with the dissemination of short videos of government affairs as the core has important practical guiding significance for guiding local grass-roots governments to build public services and public feedback.**

*Keywords—Big data statistics; short videos of government affairs; communication effect; linear regression; mainstream media*

## I. INTRODUCTION

The development of the times has promoted the progress of technology, which makes the speed of information dissemination faster and faster, and makes the content more and more [1]. In order to better catch people's attention, mass media was born with mobile short videos. Short video refers to a new video form with playback duration of less than five minutes, which can be played, shot and edited through mobile intelligent terminals, and can be shared in real time and seamlessly on social media platforms [2]. As the pace of life is getting faster and faster, people begin to pay attention to the

grasp of fragmented time. The emergence of mobile short videos meets the fragmented reading habits of the public, and also helps the new media of government affairs find a new way of political communication [3]. To improve the work of news and public opinion to a new strategic height of national governance, and in the context of the rapid development of Internet technology, the work of news and public opinion needs to pay attention to the use of new media [4]. In recent years, Party committees and governments at all levels and social groups have set up government microblogs and government official account as important government communication platforms [5]. At the same time, in view of the rise of short video and the expansion of its communication power and influence, Party committees, governments and mass organizations have settled in the short video platform, striving to create a new highland of government communication [6]. Information has become complex and difficult to distinguish between true and false. The public opinion space presents a new state of active thinking and collision of ideas. The media form carrying information and the technology of expressing information are changing with each passing day. The work of news and public opinion is facing difficult challenges and major opportunities, which poses a severe test for the government to create a good public opinion environment and grasp the requirements of Ideological and political leadership [7].

With the rapid development of short video platform represented by Tiktok, the way people obtain information has changed. More and more government agencies have settled in Tiktok short video platform. With the entertainment and light transmission characteristics of Tiktok short video platform, and giving full play to the infectious communication advantages of the combination of short video sound and painting, the dissemination of government information will be fragmented, focused and entertained [8]. Statistics is a discipline that infers and even predicts the specific situation of the measured object through the collection, collation, analysis and description of data and information [9]. Statistics is widely used in practical work, and its data collection methods and statistical analysis methods are widely used in all walks of life [10]. As a new information processing and analysis method developed with the Internet and information systems, big data also adopts certain statistical analysis methods, but it is obvious that the current big data still lacks more and more professional statistical analysis methods. In addition, big data can inspire

statistical work, and then inject some innovative thinking into statistical work, which is more conducive to the implementation of statistical work. In view of the advantages of big data statistics, this paper adopts the method of coefficient determination and linear regression in big data statistics in order to reduce the execution cost of the algorithm. Through practice, it is proved that this combination can not only reduce the calculation time, but also improve the quality and efficiency of government short video transmission optimization.

With the continuous upgrading of user demand, the short video platform has become the most popular political news public opinion field after the "two wechat ends". Among them, the "Tiktok" short video platform has the closest cooperation with government departments, and a large number of government departments have established communication positions in "Tiktok"[11]. Compared with the government news release mode in the form of pure text or graphics, short videos are more timely, more dense, more convenient to browse, more intuitive and understandable, and the social interaction experience caters to the audience's needs for information selection and self-expression. Where users gather is where the good voice of the party and government should be spread [12]. With the vigorous development of new media for government affairs, it is very urgent and necessary to adapt to the form of policy publicity in the new era, improve the efficiency of information transmission, optimize the effect of public opinion guidance, and innovate the governance mode. Relying on the authority of the official account in the social platform, it is necessary to publish authoritative news, establish a good image, and transmit positive energy through text, pictures, videos and other forms [13]. This paper establishes a feature reconstruction model for the evaluation of the dissemination effect of government short video, combs the dissemination effect and influencing factors of government short video through content analysis, analyzes the main factors by using the linear regression of big data statistics, and extracts the fuzzy feature quantity of government short video. Its innovation lies in:

*1)* This paper adopts the linear regression method in big data statistics in order to reduce the execution cost of the algorithm.

*2)* Using content analysis method, the research design is carried out according to the research paradigm of content analysis method. Referring to the influence evaluation indicators used in relevant research, this paper puts forward the parameter indicators needed to study the short video of government affairs.

This paper studies the optimization design of the dissemination effect of government short video. The architecture is as follows:

Section I is the introduction. This part mainly expounds the research background and significance of government short video communication optimization, and puts forward the research purpose, method and innovation of this paper. Section II mainly summarizes the relevant literature, summarizes its advantages and disadvantages, and puts forward the research

ideas of this paper. Section III is the method part, which focuses on the optimization design method combined with content analysis and big data statistics. Section IV is the experimental analysis. In this part, experimental verification is carried out on the data set to analyze the performance of the model. Section V, conclusions and prospects. This part mainly reviews the main contents and results of this study, summarizes the research conclusions and points out the direction of further research.

## II. RELATED WORK

Building a service-oriented government advocates simplifying administration and delegating power, and digital government affairs are increasingly showing the characteristics of convenience, humanization and intelligence, which further improves the efficiency of government services, shapes a good government image, and becomes an important channel for the government to serve the people, which is respected and accepted by government departments at all levels.

After sorting out the operation form, information release and public response of the central Tiktok account of the Communist Youth League, genton m g and sun y, according to the public's feelings about the government image of this Tiktok account, they concluded that Tiktok government affairs short video should meet the actual needs of the majority of the people, strengthen its control, and strive to improve the quality of publicity. Government affairs publicity platforms should also be diversified and use more new models [14]. Sangalli l m took the Tiktok account of the Central Committee of the Communist Youth League as an example to conduct research. The conclusion was that the type of short videos of government affairs was closely related to mass participation. Short videos of current events and hot spots could cause people to like, comment and forward. The praise and comment enthusiasm of the masses could be clearly seen in military style publicity videos. In music and emotional videos, melancholy music could promote mass comments. In terms of the title of the video, what can get more comments and forwarding from the masses is the statement [15]. With the help of the "interactive ritual chain" theory, James G M studies the interactive communication between the government Tiktok number and users. The government Tiktok number should build a perfect communication mode, and further communicate with the masses through common concerns and catering to the feelings of the masses in the process of releasing and publicizing information [16]. Okulicz kozaryn A and others believe that the interactive communication method for the future development of government Tiktok is to rely on the non popular positioning, give reasonable play to the characteristics of various platforms, launch updated output and publicity models, better disseminate information, and strengthen communication with the masses [17]. DAAS P J H et al. Took the relevant contents of the government Tiktok number of 13 central level units as the research direction, and concluded that the dissemination of department information and the construction of national image are the main contents of the central department video number, in which patriotism is widely disseminated. At the same time, these government affairs Tiktok numbers are also good at grasping the main time points and important events with high popularity, and arousing the

resonance of the masses through videos. These measures have increased the attention of video numbers and brought new enlightenment to government affairs communication [18]. Through the case study of the innovative characteristics of "Shenzhen Energy", Dozier J studied the content and technological innovation of government Tiktok short video, built a government characteristic culture, made government services deeply rooted in the hearts of the people, and realized the sustainable development of government Tiktok short video [19]. In the context of the outbreak of the COVID-19, Scanlon D P and others studied the government Tiktok number in combination with the particularity of the epidemic environment. After analyzing the content and function of the government Tiktok number, they proposed that the government Tiktok number was not beneficial to the control of public opinion and guidance of events, and the relevant government work members should pay attention to and solve such problems [20]. Zhang y et al. analyzed from a new perspective, that is, the new communication characteristics emerging in the context of media integration. The "decentralization" and "de organization" characteristics of new media have enhanced the activity of users and the transparency of the social public opinion environment. Under this background, the new media of government affairs is the social product of actively adapting to and following the laws of the Internet and actively innovating service methods [21]. Dunson, David B pointed out that in the context of media convergence, media operators should have Internet thinking, accept and accommodate new media with an open mind, and expounded the significance of building a new media matrix for government affairs. On this basis, he pointed out some problems that still exist in building a new media for government affairs [22].

In the dissemination process of government short videos, user feedback is crucial. However, existing research lacks a deep understanding of the emotional tendencies and meanings of user comments, likes, shares, and other behaviors, and fails to fully grasp the true attitudes and feelings of users towards government short videos. At present, research mainly focuses on the cultural background of China, and there is relatively little research on the dissemination effects of government short videos in other countries and regions. Therefore, there are obvious shortcomings in cross-cultural comparison. In addition, there is a lack of sufficient research on the comparison of the dissemination effects of different short video platforms. Although some studies have proposed strategies to improve the effectiveness of government short video dissemination, these strategies often lack practical application value and fail to effectively translate into specific operational suggestions or solutions. Therefore, how to translate research results into practical applications is still an urgent problem to be solved.

In order to make up for these shortcomings, future research needs to further expand the depth and breadth of data processing, explore user feedback in depth, strengthen cross-cultural and cross platform comparative research, and improve the practical application value of improvement strategies. Through these efforts, we can comprehensively and accurately evaluate the dissemination effect of government short videos, and provide more valuable suggestions for practical

applications. In the face of the blowout of public opinion, the traditional government affairs communication mode, which used to use the media to speak out, was stretched out and unsustainable, falling into the communication dilemma of "talking to yourself" and the decline of the government's credibility, which led to the failure of the guidance of government public opinion and the "failure" of government affairs communication. This paper proposes an optimization design scheme for the evaluation and promotion research of the dissemination effect of government short video based on big data statistics. By using the methods of content analysis and data analysis, this paper explores the significant influencing factors that affect the dissemination of government short video content, and then obtains the dissemination strategies to improve the dissemination of government short video content through data conclusions. Optimize the short video of government affairs from the aspects of video content and editing techniques, so as to make the planning of short video of government affairs more reasonable.

## III. METHODOLOGY

### A. Classify and Quantify Through Content Analysis to Analyze the Dissemination of Government Short Videos

In the online world, users are fully exposed to a large amount of information. From the sending of information to the acceptance of information, they will be affected by a variety of factors [23]. In practical applications, the extraction and analysis of unstructured and multidimensional big data have a wide range of application scenarios. In market analysis, market trends and consumer demand can be predicted by analyzing consumer online behavior data. In urban planning, multidimensional big data analysis can be used to evaluate the development status of cities and provide scientific basis for policy formulation [24]. With the popularization of digital media, government agencies have also begun to use short video platforms to interact and disseminate information with the public. In order to better evaluate the dissemination effect of government short videos, big data statistics have become an important tool [25]. Relying solely on big data statistics is not enough. When evaluating the effectiveness of government short video dissemination, multiple dimensions need to be considered. For example, in addition to basic viewing data, the theme, style, and audience characteristics of video content can also be analyzed. By understanding the interests and needs of the audience, the content and quality of short videos can be further optimized [26].

While specific media content only appears once, users are not forced to pay attention to this information as experimental participants, and will not be hinted by psychology. On the contrary, their feedback data for information content only depends on the role and results of their personal psychological and social characteristics. This paper boldly believes that the Internet world is like a huge natural laboratory, with interference from various factors that will affect the communication process. However, because its subjects are netizens, the base number is very large. In such a large experiment, many subtle factors can be ignored, and the real user data is an embodiment of the effect of network information communication, and all netizens show their

attitude towards information. Taking real user data as the indicator of communication effect, although user data is only the embodiment of information acting on cognition and attitude, and has not yet penetrated into the user's behavior level, it cannot be denied that high viewing volume, high praise number, high comment number and high forwarding number have achieved the communication purpose to a certain extent, and taking into account the practical operability of the research, the measurement of its explicit data is appropriate and reasonable. Fig. 1 shows the proposed content dissemination capability model.



Fig. 1. Suggested Model of Content Communication Power

This paper refers to the classification of government short videos in some official accounts of government accounts, and classifies government short videos according to three levels: administrative level, industry system and content nature. It is divided into ministries and commissions, provincial, municipal and county levels according to different administrative levels; According to the industry system, it is mainly divided into public security, fire protection, procuratorate, court, Communist Youth League, financial media, cultural tourism, etc; According to the nature of the content, it can be divided into science popularization, publicity, interaction, news and story as shown in Table I.

TABLE I. CLASSIFICATION AND INDUCTION OF SHORT VIDEOS ON GOVERNMENT AFFAIRS

| Government short video classification | |
|---|---|
| Administrative level | Ministerial level |
| | Provincial level |
| | Municipal level |
| | County-level |
| Administrative system | Political and Legal Committee |
| | Public Security |
| | Procuratorial class |
| | Court class |
| | Judicial category |
| | Communist Youth League |
| | Financial media |
| | Health |
| | Cultural tourism |
| | Women's Federation |
| Content nature | Popular science |
| | Publicity |
| | Interactive class |
| | News |
| | Stories |

Fig. 2.   Model of factors affecting the dissemination of government short video content.

As a new form of government information dissemination, government short video widens the space and channels of government short video information dissemination. In terms of functional positioning, compared with China's traditional government short video new media, its main function is positioned as "government online performance ability" and "online government at the fingertips". Its core function is information disclosure and dissemination, rather than government services.

Combined with the characteristics of short video of government affairs and the attributes of short video platform, this study extracts the content theme, personas, video type, video emotion, patriotic emotion expression, video duration, use of background music, subtitles with dialogue, subtitles with special effects, symbols used in titles, network terms used in titles, and pragmatic expressions of titles in the production of short video content of government affairs in content planning, post editing, and operation release, Video release time and personalized communication means are 14 influencing factors, and a model of influencing factors of the dissemination of short video content of government affairs is constructed. The model is shown in Fig. 2.

### B. Optimize the Communication Effect of Government Short Video Based on Big Data Statistics

The rapid development of new media in the era of big data has made a great change in the traditional mode of information transmission, and gradually affected people's digital lifestyle and the habit of contacting the mass media. Since the major short video platforms have entered people's lives, short video software has become a very important part of the mobile video industry. The rapid development of the short video industry has gradually enriched people's lives. Compared with traditional words and pictures, short video, a new media form, is more vivid, intuitive and entertaining, which is consistent with the behavior of people sharing and participating in information dissemination on the network platform in the current era of big data.

Determination coefficient

$R^2$ is also called the determination coefficient of the equation, which indicates the interpretation degree of variables $X$ to $Y$ in the equation. The value of $R^2$ is between [0,1], and the closer $R^2$ is to 1, the stronger the explanatory ability of $X$ to $Y$ in the equation. Usually, $R^2$ times 100% is used to express the percentage of change in the interpretation $Y$ of the regression equation.

Taking the simplest univariate linear regression analysis as an example, this paper expounds the basic principle of the determination coefficient.

As above, the observation data are:

$$(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n), \tag{1}$$

The determination coefficient of the unary linear problem is desired.

$$y_i = a + bx_i + \varepsilon_i = \hat{y}_i + \varepsilon_i \tag{2}$$

where, $\hat{y}_i$ is the calculated quantity, $i = 1,2,\ldots,n$. Based on the average value $\overline{y}$ of the explained variable $y_i$, the above formula can be transformed into,

$$y_i - \hat{y} = (\hat{y}_i - \overline{y}) + \varepsilon_i = (\hat{y}_i - \overline{y}) + (y_i - \overline{y_i}) \tag{3}$$

Using the SRF sample regression function, there are,

$$\sum_{i=1}(y_i - \overline{y})^2 = \sum_{=1}(\hat{y}_i - \overline{y})^2 + \sum_{=1}(y_i - \overline{y_i})^2 \tag{4}$$

Among them, $\sum_{i=1}^{n}(y_i - \overline{y})^2$ is called the sum of squares of total variables or total force differences, $\sum_{i=1}^{n}(\hat{y}_i - \overline{y})^2$ is called the sum of squares of regression, and $\sum_{i=1}^{n}(y_i - \hat{y}_i)^2$ is called the sum of squares of residuals.

Because the fitting effect of the sample regression line on the observed value depends on the distance between the sample observation value and the regression line, that is, the proportion of the sum of the squares of the regression in the sum of the squares of the total deviations. Therefore, the judgment coefficient can be obtained.

$$R^2 = \frac{\sum_{i=1}^{n}(\hat{y}_i - \overline{y})^2}{\sum_{i=1}^{n}(y_i - \overline{y})^2} = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \overline{y})^2} \tag{5}$$

Pass $t$ inspection

From the above, it is easy to prove that $Var(\varepsilon_i|X) = E(\varepsilon_i^2) = \sigma^2$ and $\hat{a}, \hat{b}$ are the designs of $a$ and $b$

$$Var(\hat{b}) = \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \overline{x})^2} \tag{6}$$

$$Var(\hat{a}) = \sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \overline{x})^2} \tag{7}$$

And $\hat{a}, \hat{b}$ obey normal distribution respectively

$$\hat{a} \sim N(a, \sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \overline{x})^2}) \tag{8}$$

$$\hat{b} \sim N(b \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \overline{x})^2}) \tag{9}$$

Normalize normal random variables $\hat{a}$ and $\hat{b}$ to obtain

$$z_1 = \frac{\hat{a} - a}{SE(\hat{a})} = \frac{\hat{a} - a}{\sqrt{\sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \overline{x})^2}}} \tag{10}$$

$$z_2 = \frac{\hat{b} - b}{SE(\hat{b})} = \frac{\hat{b} - b}{\sqrt{\frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \overline{x})^2}}} \tag{11}$$

$SE$ represents the standard deviation of the variable. At this time, $z_1$ and $z_2$ obey the standard normal distribution.

However, the variance $\sigma^2$ of random disturbance term $\varepsilon$ is unknown, and it can only be estimated unbiased with $\hat{\sigma^2} = \sum_{i=1}^{n} \varepsilon_i^2 / (n-2)$, available at this time.

$$SE(\hat{a}) = \sqrt{\sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n} (x_i - \overline{x})^2}} \tag{12}$$

$$SE(\hat{b}) = \sqrt{\sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n} (x_i - \overline{x})^2}} \tag{13}$$

In the case of small samples, it is easy to prove that $(\hat{a} - a)/SE(\hat{a})$ and $(\hat{b} - b)/SE(\hat{b})$ no longer obey the standard normal distribution, but obey the $t$ distribution with a degree of freedom of $n - 2$, that is

$$t = \frac{\hat{a} - a}{SE(\hat{a})} = \frac{\hat{a} - a}{\sqrt{\sigma^2 \frac{\sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n} (x_i - \overline{x})^2}}} \sim t(n-2) \tag{14}$$

$$t = \frac{\hat{b} - b}{SE(\hat{b})} = \frac{\hat{b} - b}{\sqrt{\frac{\sigma^2}{\sum_{i=1}^{n} (x_i - \overline{x})^2}}} \sim t(n-2) \tag{15}$$

## IV. RESULT ANALYSIS AND DISCUSSION

Whether it is the former text era, the glorious newspaper era, and today's new media era with the emergence of science and technology, content is the prerequisite and primary concern of the media, but also its necessary standard. To succeed in short video, it is necessary to put the production of content in the first place. Only in this way can we get the audience's love.

As can be seen from Table II, the average value of the four options is >3, indicating that the audience generally likes the short video of government affairs, which is a reason to attract attention from both the professionalism of the content and the style type.

For the "online celebrity" in the short video of government affairs, the video number must rely on high-quality content and different forms of communication to enhance its competitiveness. High quality content is the top priority among them. The audience gets the hot news they are interested in or the knowledge they need to learn from short videos. In this era of complex information, people's fast-paced life is prone to anxiety, which leads to being hoodwinked by some new media that convey wrong values. The short videos of the Central Committee of the Communist Youth League help the audience understand the whole story by interpreting the recent hot information. It is conducive to improving the social cognition of the audience, and these contents can really meet the needs of the audience and get the audience's love.

As a social platform, the most remarkable feature of short video platform is that it can quickly realize the function of one click forwarding in a short time, and interpersonal networking naturally plays a role in promoting the dissemination of content. It is an effective way to maintain interpersonal communication to transmit information that is fun and interesting in real life or useful to friends and can produce practical effects. Therefore, the government affairs short video should enlarge its own social communication elements to improve the audience's desire to share, so that the audience can analyze and discuss this phenomenon and improve its communication effect and influence. The importance of the audience's demand for social interaction can be seen from Table III. In recent years, domestic mobile social networking platforms have developed rapidly, and mainstream media and communication media have entered mobile social networking platforms.

TABLE II. CONTENT POPULARITY STATISTICS

| | N | Minimum value | Maximum value | Mean value | Standard deviation |
|---|---|---|---|---|---|
| Rich and interesting content | 375 | 3 | 5 | 3.72 | 0.745 |
| Content of military image displayC | 375 | 2 | 5 | 3.66 | 0.777 |
| Content about the deeds of model figures | 375 | 3 | 5 | 3.74 | 0.764 |
| Content of transmitting positive energy | 375 | 3 | 5 | 3.65 | 0.761 |
| Deliver the content of excellent traditional culture | 375 | 2 | 5 | 3.71 | 0.752 |
| Valid N | 375 | | | | |

TABLE III. ANALYSIS OF SOCIAL INTERACTION FACTORS

| | N | Minimum value | Maximum value | Mean value | Standard deviation |
|---|---|---|---|---|---|
| Interact with netizens in the comment area | 375 | 1 | 5 | 3.86 | 0.756 |
| More interested in the content of the popular list | 375 | 2 | 5 | 3.75 | 0.743 |
| More interested in sharing likes with friends | 375 | 3 | 5 | 3.77 | 0.774 |
| Share your views and meet social needs | 375 | 2 | 5 | 3.82 | 0.759 |
| Valid N | 375 | | | | |

Short videos create a platform for the audience to "speak freely" through star attraction and interesting videos. Government related institutions also see the significant advantages of short videos, which can trigger users' enthusiasm for reprinting, liking and commenting while spreading. It is very consistent with the platform attribute of current social media, and realizes the transmission of positive energy in the platform that seems to be threatening entertainment, You can not only get some useful information, but also experience relaxation and entertainment. This also shows that in the current era of integrated media, people's offline activities are gradually reduced, and online social relations have an important impact on people. In today's diversified and complex information, group sharing and interaction are not simply about content, but also meet people's desire to share in personalized

communication. Social platforms need to strengthen their social attributes, so that people can more easily deliver the content they are interested in, promote multi-level content dissemination, make good content benefit more people, expand influence, and achieve good communication results.

In order to have a targeted and effective understanding of the use of government short video, and provide an effective reference for the development status and future development strategies of government short video, the influencing factors were investigated. The specific distribution of six aspects: the usefulness of the information, the interest of the information provided, the comprehensibility of the information provided, the speed of information update, the beauty of the overall video style, and whether the government Tiktok has played its value is shown in Fig. 3.



Fig. 3. Use value of government short video.

Based on all the data, it can be concluded that most of the respondents are passive in contacting the government Tiktok number, and they do not take the initiative to understand the situation of concern. Even many people do not know the existence of the government Tiktok number, or their cognition of the government Tiktok number is vague, and their demand and desire for the government Tiktok number is not very high, let alone the situation and desire of interactive participation. There are not many people who want to pay attention to the government Tiktok number, for the understanding of the government Tiktok number, we hope to encounter information or promote the page. If you need to pay attention, the areas of attention tend to be relevant to yourself.

Through the research samples obtained above, we conduct empirical research on the influencing factors of the content dissemination of government short videos, and summarize the main factors that significantly affect the content dissemination of government short videos: content theme, personas, video type, video emotion, a total of four influencing factors. This chapter will discuss and analyze the specific research data and

related theories, and then summarize the research on improving the dissemination of government short video content.

It is found that the number of likes, comments and forwarding of short government videos with content themes of remembering history is significantly higher than that of other content themes; the forwarding number of short government videos with working dynamic content topics is significantly lower than that of other content topics.

According to the results of stepwise linear regression in Fig. 4, the content theme of remembering history will significantly and positively affect the number of likes, comments and forwards; the content subject of working dynamic class will significantly negatively affect the number of forwarding. This significant difference can be explained by using and satisfaction theory. As active rational individuals, current short video users have the right to choose. Users can selectively watch, like, comment and forward short video messages according to their personal needs, and the feedback data of the short video message corresponds to the individual's satisfaction with their needs to a certain extent.

Fig. 4. Correspondence between content likes, comments and average forwarding data.



Fig. 5. Correspondence between characters and the average value of likes, comments and forwarding data.

According to the results of stepwise linear regression in Fig. 5 above, short videos of government affairs without people will significantly and positively affect the number of likes and forwards; if the persona is a "specific person", it will significantly negatively affect the number of forwards.

Government affairs short videos without people are generally macro scene descriptions of major events or text event notifications. The personas in the short video represent the narrative perspective, which is generally the perspective of the narrator's storytelling, which also contains the emotional tendency hidden by the narrator. The more such personas appear in the short video of government affairs, it also shapes and promotes the quality spirit of such personas, but the short video content with personas is more didactic, whether it is the third person perspective of telling persona stories, It is also the first person perspective of the personas' own "story telling",

and the short videos with specific personas have a stronger meaning of "value leading". When there is no persona in the short video of government affairs, the user watches it from a subjective perspective, so that the user can be immersive and empathic, and obtain the cognition of the information from the aspects of vision, hearing and even spiritual feeling.

According to the results of stepwise linear regression in Fig. 6, video emotion as "moving" emotion will significantly and positively affect the number of likes. The epidemic situation in the century and the changes in the century are intertwined, and severe challenges and major difficulties coexist. However, the indomitable people overcome the difficulties together, reflecting the national speed, demonstrating the national strength, and creating national miracles.

Fig. 6. Correspondence between video feelings and the average value of likes, comments and forwarded data.



Fig. 7. Correspondence between publishing time and average value of likes, comments and forwarded data.

In Fig. 7, according to the results of stepwise linear regression, the release of short government videos at 8-10 will significantly and positively affect the number of likes and forwards. This is because the target users' active time is different, and the release time of a single short video content is significantly different from the data of content dissemination. Theoretically, the target user group has the same information receiving habits, and the active time of receiving information is also relatively consistent. The video content released during the peak period of user activity is relatively more likely to become popular. According to the report on the release time of Tiktok short visual frequency released by Kasi data, it is pointed out that the videos released by "Tiktok online celebrity" generally at 17-18 o'clock are easier to gain interaction, and the videos from 11-12 o'clock at noon are also good. The result of this paper is that the government short video released at 8-10 a.m. will get more likes and forwards, which means that the active time of the target user group of the government short video is 8-10 a.m.

## V. CONCLUSIONS

This paper proposes an optimization design scheme for the evaluation and improvement of the dissemination effect of government short video based on big data statistics. The basic situation of government video is obtained through content analysis, and then the judgment coefficient and linear regression in big data statistics are used to extract common factors to improve the dissemination effect of government short video, so as to improve the dissemination influence of government short video. Finally, simulation test and analysis are carried out. Simulation results show that the proposed algorithm has certain accuracy, which is 8.24% higher than the traditional algorithm. This result fully shows that from the traditional media era to the current Internet new media era, user attention has become an important resource, and it is undeniable that high-quality content has always been a magic weapon to attract user attention. "Content is king" is not out of date, but puts forward higher requirements. How to carry content, express content, and disseminate content has also become the content itself. This study systematically takes the

dissimilated content as the foundation, the high-quality system as the guarantee, and the linkage communication as the advantage from the three aspects of content, content production and operation, in order to stand out from the redundant information containment in this information age with the rise of we media and social media, seize the attention of users, and improve the content communication power of government short videos. As the short video of government affairs is a new product of short video and a new content carrier of government affairs communication, its academic research and practical development are still in the initial stage. How to carry out effective government affairs communication of short video of government affairs is a new problem faced by both academia and industry. The data of this study mainly comes from various short video platforms, such as Tiktok, Kwai, etc. Although these platforms have a large user base in China, there may still be specific groups or regions using other platforms, which may affect the comprehensiveness of our data. In big data analysis, the quality and accuracy of data are key issues. Although we have employed various methods to ensure the accuracy and completeness of the data, there may still be some errors or omissions that may have an impact on the research results.

In the future, in addition to mainstream short video platforms, data from other platforms or social media can also be considered for more comprehensive analysis by adopting more advanced data cleaning and preprocessing techniques which will help to improve the accuracy and completeness of data.

## COMPETING OF INTERESTS

The authors declare no competing of interests.

## AUTHORSHIP CONTRIBUTION STATEMENT

Man Xu: Writing-Original draft preparation, Conceptualization, Supervision, Project administration.

## DATA AVAILABILITY

On Request

## DECLARATIONS

Not applicable

## REFERENCES

[1] S. Hong et al., "Constraining cosmology with big data statistics of cosmological graphs," Mon Not R Astron Soc, vol. 493, no. 4, pp. 5972–5986, 2020.

[2] S. McGrath, X. Zhao, Z. Z. Qin, R. Steele, and A. Benedetti, "One-sample aggregate data meta-analysis of medians," Stat Med, vol. 38, no. 6, pp. 969–984, 2019.

[3] C. G. Rossa, "The effect of fuel moisture content on the spread rate of forest fires in the absence of wind or slope," Int J Wildland Fire, vol. 26, no. 1, pp. 24–31, 2017.

[4] J. F. David and S. A. Iyaniwura, "Effect of Human Mobility on the Spatial Spread of Airborne Diseases: An Epidemic Model with Indirect Transmission," Bull Math Biol, vol. 84, no. 6, p. 63, 2022.

[5] B. W. Pitcher and A. J. R. Kent, "Statistics and segmentation: using big data to assess Cascades arc compositional variability," Geochim Cosmochim Acta, vol. 265, pp. 443–467, 2019.

[6] A. O. Afolayan, J. S. Mandeep, M. Abdullah, and S. M. Buhari, "Statistics of spread F characteristics across different sectors and IRI 2016 prediction," Advances in Space Research, vol. 64, no. 10, pp. 2154–2163, 2019.

[7] F. Pimont, J. Ruffault, N. K. Martin-StPaul, and J.-L. Dupuy, "Why is the effect of live fuel moisture content on fire rate of spread underestimated in field experiments in shrublands?," Int J Wildland Fire, vol. 28, no. 2, pp. 127–137, 2019.

[8] A. Grzybowski and M. Mianowany, "Statistics in ophthalmology revisited: the (effect) size matters," Acta Ophthalmol, vol. 96, no. 7, pp. e885–e888, 2018.

[9] Z. Guo, M. Gully-Santiago, and G. J. Herczeg, "The Effect of Spots on the Luminosity Spread of the Pleiades," Astrophys J, vol. 868, no. 2, p. 143, 2018.

[10] M. G. Cruz, A. L. Sullivan, R. Bessell, and J. S. Gould, "The effect of ignition protocol on the spread rate of grass fires: a comment on the conclusions of Sutherland et al.(2020)," Int J Wildland Fire, vol. 29, no. 12, pp. 1133–1138, 2020.

[11] Y. Yuan et al., "Key frame extraction based on global motion statistics for team-sport videos," Multimed Syst, vol. 28, no. 2, pp. 387–401, 2022.

[12] W. P. Nobis et al., "The effect of seizure spread to the amygdala on respiration and onset of ictal central apnea," J Neurosurg, vol. 132, no. 5, pp. 1313–1323, 2019.

[13] R. Cao, "Comments on: Data science, big data and statistics," Test, vol. 28, no. 3, pp. 664–670, 2019.

[14] M. G. Genton and Y. Sun, "Comments on: Data science, big data and statistics," Test, vol. 28, no. 2, pp. 338–341, 2019.

[15] L. M. Sangalli, "The role of statistics in the era of big data," Stat Probab Lett, vol. 136, pp. 1–3, 2018.

[16] G. M. James, "Statistics within business in the era of big data," Stat Probab Lett, vol. 136, pp. 155–159, 2018.

[17] A. Okulicz-Kozaryn and J. M. Mazelis, "More unequal in income, more unequal in wellbeing," Soc Indic Res, vol. 132, no. 3, pp. 953–975, 2017.

[18] P. J. H. Daas, M. J. Puts, B. Buelens, and P. A. M. van den Hurk, "Big data as a source for official statistics," J Off Stat, vol. 31, no. 2, pp. 249–262, 2015.

[19] J. Dozier, "Revisiting topographic horizons in the era of big data and parallel computing," IEEE Geoscience and Remote Sensing Letters, vol. 19, pp. 1–5, 2021.

[20] D. P. Scanlon and M. B. Stephens, "Tests, surgical masks, hospital beds, and ventilators: Add big data to the list of tools to fight the coronavirus that are in short supply," Am J Manag Care, vol. 26, no. 6, pp. 241–244, 2020.

[21] Y. Zhang et al., "Mobile social big data: Wechat moments dataset, network applications, and opportunities," IEEE Netw, vol. 32, no. 3, pp. 146–153, 2018.

[22] D. B. Dunson, "Statistics in the big data era: Failures of the machine," Stat Probab Lett, vol. 136, pp. 4–9, 2018.

[23] L. Chen and Y. Zhou, "Quantile regression in big data: A divide and conquer based strategy," Comput Stat Data Anal, vol. 144, p. 106892, 2020.

[24] Adnan, K., & Akbar, R. (2019). An analytical study of information extraction from unstructured and multidimensional big data. Journal of Big Data, 6(1), 1-38.

[25] Soomro, K., Bhutta, M. N. M., Khan, Z., & Tahir, M. A. (2019). Smart city big data analytics: An advanced review. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 9(5), e1319.

[26] Rehman, A., Naz, S., & Razzak, I. (2022). Leveraging big data analytics in healthcare enhancement: trends, challenges and opportunities. Multimedia Systems, 28(4), 1339-1371.

# Improved Ant Colony Algorithm Based on Binarization in Computer Text Recognition

Zhen Li

School of Forensic Science and Technology, Criminal Investigation Police University of China, Shenyang, 110854, China

*Abstract*—Pheromones, path selection, and probability transfer functions are the main factors that affect the performance of computer text recognition. The path selection function is the most important factor affecting the recognition rate. In response to the difficulties in path selection and slow algorithm convergence in the text recognition, an edge detection algorithm based on improved ant colony optimization algorithm is proposed. The strong denoising performance of the ant colony optimization algorithm reduces the interference of textured backgrounds. The edge extraction effect is analyzed in the connected domain to overcome complex effects. Finally, the improved Otsu binarization algorithm is used to recognize the text. According to the results, the proposed method could effectively preserve the edge information of characters in images. The positioning effect of the text area was good. The accuracy rate reached around 85%. The tuned threshold improved the binarization effect. The text recognition rate of the improved ant colony algorithm proposed in the research has generally reached 80%, with good text positioning accuracy and recognition rate, which has great practical significance in computer text recognition.

*Keywords*—*Binarization; ant colony algorithm; text recognition; edge detection; Otsu algorithm*

## I. INTRODUCTION

The core of computer text recognition is to recognize characters. Therefore, binarization of character images is the prerequisite and foundation for it [1-2]. The high degree of deformation and discontinuity in text images makes binarization techniques in image processing very difficult. Most existing character recognition algorithms analyze character images to determine the corresponding character category for each character. The main steps include image processing, feature extraction, character segmentation, and character recognition. The feature extraction part is the most important part of the entire algorithm. Its main task is to transform the text image into a series of feature vectors for text recognition based on the analysis of the text image. In the text recognition, the image binarization quality directly affects the quality of text recognition result. Therefore, how to achieve binarization is a key link in the entire text recognition process [3]. The existing binarization algorithms include maximum inter class variance method, threshold segmentation method, and threshold division method. These algorithms have certain limitations in practical applications. They cannot meet the binarization requirements of text and images in complex backgrounds [4-6]. Due to factors such as pheromone concentration, path selection, and probability transfer function, there are redundant pheromones and invalid paths in

text images before binarization, resulting in noisy or blurry areas in the binarized text image [7-8]. To obtain high-quality binarization results, an adaptive method must be used to smooth the binarized image. However, the smoothing methods used in existing algorithms to some extent increase the noise points or blurry areas of text images. In the multimedia information retrieval, the image and video search in the search engine is still based on keywords, and the image is manually annotated first. When the user enters the keyword to search, in fact, the search engine only maps the results to the corresponding pictures. However, for complex and changeable video and various kinds of text, these methods have limitations, and are difficult to achieve a practical application level. In order to extract text in multimedia quickly and effectively, a method of "image search" is proposed, that is, by extracting multi-dimensional information such as texture, color and shape of the target image, the most similar image from the image library is find according to a specific algorithm. Texts are extracted from the video stream to realize content-based video retrieval.

The improved ant colony algorithm (ACA) is combined with the edge detection algorithm to solve the problem of slow convergence and local optimal solution. In the binarization stage, an improved Otsu binarization algorithm is proposed on the basis of the traditional binarization algorithm, and the corresponding binarization threshold is obtained by using the traditional Otsu algorithm. The threshold is fine-tuned on the basis of the traditional algorithm to improve the binarization effect.

The article conducts research from six sections. Related works is given in Section II. Section III is a review of research on improved ACA in computer text recognition. Section IV constructs the text recognition method based on improved ACA. Section V is to verify the performance of the proposed method. Section VI is the conclusion.

## II. RELATED WORKS

ACA is extensively applied in combinatorial optimization problems, such as travel salesman problems, vehicle routing problems, graph coloring problems, and network routing problems. To promote economic development, many scholars have conducted research on this. Yi et al. proposed an improved ACA for task scheduling problems in information physics systems. The adaptive and mutation strategies were adopted to reduce solution time and accelerate the convergence of information physics systems. After numerical simulation, the improved algorithm could effectively improve the local optimization ability. It had good adaptability and

stability [9]. To solve the long time consumption and multiple intermediate nodes in traditional path guidance methods in path planning results, Tang et al. optimized the potential field of ant colony from three aspects: potential field function, pheromone update process, and heuristic function. The ability to avoid obstacles and optimize was enhanced. The results showed that this method had fewer midpoints and the shortest path on each route [10]. Zhu et al. fused the artificial potential field method with the ant algorithm to overcome the slow convergence speed and local optimal problem of the ACA. The induction heuristic factor was used to dynamically adjust the state transition law, resulting in higher global search ability and faster convergence speed of the algorithm. After verification, this method could effectively determine whether a collision is occurred and take obstacle avoidance measures [11]. Yu et al. introduced a new heuristic clustering algorithm called co-evolutionary chain to improve the accuracy and stability of ACA. Ant colony clusters were divided for the balance of convergence and speed. The combination of group co-evolution and link dimension reduction improved the diversity and stability, separating it from the local optimal problem [12]. Hwang et al. proposed an ACA bidirectional Long Short Term Memory (LSTM) network model. The existing physiological signal information was fully utilized. The ACA was applied to find the optimal emotion recognition features, improving the emotion recognition ability of existing LSTM cell states. The final results indicated that the model had good valence performance [13].

Text recognition is widely used in daily life. Currently, many applications have been put into practice, such as documents conversion, license plate recognition, photo search, target translation, etc. Guptha et al. proposed a new deep learning based automatic character recognition model to address the difficulty of handwritten character recognition. Gaussian filtering and tilt detection techniques were used for preprocessing handwritten images. Projection contour and threshold segmentation techniques were used to segment individual lines and characters from denoised images. Finally, characters were classified using the LSTM [14]. Liu et al. developed a recognition algorithm to address the low character recognition efficiency and accuracy. A feature model was developed by combining histogram Gabor features with grid level features. Then, a deep belief network was applied to train the feature model. Finally, a probability model was used to judge the recognition symbols. After verification, it had higher accuracy and better performance [15]. According to the image information obtained by a 3D camera, Alam et al. proposed a new method for character recognition based on finger joint tracking system. The distance between the tip of the thumb and the joint of another finger was used for calculation. The Euclidean distance threshold and geometric slope technology were used to recognize numbers, letters, characters, special keys, and symbols. After verification, the overall recognition accuracy was over 90%. The recognition time for each character was less than 60 milliseconds [16]. In response to the low accuracy of character recognition in natural scenes, Chandio et al. extracted image features by cutting characters. Then, the obtained features were transmitted to the machine learning classifier for classification and character recognition through directed gradient histograms. The accuracy of this

method reached 78.52% [17]. Lee et al. proposed a real-time character recognition algorithm. Based on the architecture of an improved local binary mode shallow depth convolutional neural network (CNN), it combined the manual feature preprocessing and character learning in CNN supervised advanced functions. Networks with different depths were applied for learning. The learned features were used for classification. The algorithm had good performance [18].

In summary, researchers have proposed many methods to improve the accuracy of computer text recognition. They have also achieved certain results. However, accuracy and efficiency are still lacking. Therefore, by improving the ACA based on binarization, it is expected to quickly and effectively increase the accuracy of text recognition.

## III. TEXT RECOGNITION METHOD BASED ON IMPROVED ACA

A text recognition method based on improved ACA is proposed, which optimizes the edge detection algorithm and improves the accuracy of character region localization in images. Then, the Otsu algorithm is improved to find the optimal binarization threshold and enhance the binarization effect.

### A. Character Localization Based on ACA

ACA is a simulation optimization method that mimics ants' foraging behavior. It has been widely applied in practical applications, especially in combination optimization, which has achieved great achievements. This algorithm utilizes pheromones to remember the best route, thereby strengthening the route and finding the best one. The study transforms the problem of text recognition into solving the optimal path selection problem. The improved ACA is applied to solve linear optimization problems. The characteristics of the optimal solution are determined based on ACA. The movement path of ants between pixels in the image is used to represent the edges of the image. When an ant moves from one pixel to another, it leaves pheromones along the way (the release rules of pheromones include a series of information such as image gradient and image color). These heuristic information are used to determine pixel regions with obvious edge features between the paths passed.

Edge detection has important implications in image processing. Its results directly affect image detection and content recognition [19]. ACA is used for edge detection of text images. The movement path of ants between pixels is used to represent the edges of the image. These heuristic information can determine the pixel regions with obvious edge features between the paths passed. When the initial value of pheromone concentration $\tau_{i,j}$ is not 0 and it moves, it will volatilize with a volatilization rate $\rho$ over time. The relationship is shown in Eq. (1).

$$\tau_{i,j}(t+1) = (1-\rho)\tau_{i,j}(t) + \sum_{k-1}^{m} \Delta\tau_{i,j}^{k} \tag{1}$$

In Eq. (1), $\rho \in (0,1]$. $\Delta\tau_{i,j}^{k}$ represents the amount of pheromones added during ant movement, as shown in Eq. (2).

$$\Delta \tau_{i,j}^{k} = \begin{cases} \dfrac{\eta_{i,j}}{255}, k \in (i,j) \cup \eta_{i,j} > T \\ 0, \ \text{Or else} \end{cases} \tag{2}$$

In Eq. (2), $\eta$ represents heuristic information. $T$ represents the threshold of binarization, which affects pheromone map updates and detection. The transfer probability for ant movement is expressed in Eq. (3).

$$p_{(i_0,j_0)(i,j)}^{n}(i,j) = \frac{\left(\tau_{i,j}^{n-1}\right)^{\alpha}\left(\eta_{i,j}\right)^{\beta}}{\sum_{(i,j \in allowed_k)}\left(\left(\tau_{i,j}^{n-1}\right)^{\alpha}\left(\eta_{i,j}\right)^{\beta}\right)} \tag{3}$$

In Eq. (3), $(i,j)$ represents a pixel. $\alpha$ represents the pheromone heuristic factor. A higher value indicates a greater likelihood of repeated searches. $\beta$ stands for the expected heuristic factor. The numerical value determines the probability of ants choosing the local shortest path. $allowed_k$ represents the pixels that ants can pass through.

The implementation process of the edge detection ACA is displayed in Fig. 1. Firstly, initialization is performed. On this basis, each ant generates heuristic information based on the pheromone values of adjacent pixel points. At the same time, the transition probability is calculated. The pixel points that the next path will experience are determined based on the transition probability. The pheromone value is updated after entering a new pixel. When the ant colony moves to a new pixel position, it can update the overall pheromone map. When searching for paths, taboo table permissions are introduced to avoid ant colonies from repeatedly reaching the same pixels on the same path. If the next pixel entered by the ant colony is empty, it is randomly assigned to a pixel data to start over. Otherwise, the algorithm ends and the input grayscale image is converted into a binary image.

In the edge detection, ACA always has the slow convergence speed and the problem of local optimal solutions. The main difficulties in algorithm convergence are as follows. The distribution of initial points is random, without considering initial pheromones, resulting in slow search speed, high computational complexity, and long time consumption [20]. If it is trapped in a local optimal solution, it will cause the search process to pause and make it difficult to obtain a global optimal solution. Traditional ACAs do not limit the initial position of ants. There is no restriction on the initial layout of ants, but a random selection method is used. Therefore, ants continuously iterate during the search process, ultimately obtaining the optimal path [21]. However, using a random distribution method can lead to a large number of ants blindly searching under uncertain boundaries, thereby reducing search efficiency and lacking targeted edge detection. In view of this, a new 5×5 neighborhood based ACA is proposed to address the problems in ACA. This algorithm can better reflect the true boundary characteristics and effectively achieve the optimal localization of ant colonies.

The method for obtaining gradient information is as follows. The pixel points in the $allowed_k$ image are set to $(i,j)$. Their grayscale value is $f(i,j)$. Then the grayscale gradient $\psi_T(i,j)$ of the pixel can be expressed as Eq. (4).

$$\psi_T(i,j) = \frac{\sqrt{\left[f(i+1,j+1)-f(i-1,j-1)\right]^2 + \left[f(i+1,j-1)-f(i-1,j+1)\right]^2}}{255} \tag{4}$$

The gradient information of the image is used to extract pixels with significant differences. Meanwhile, based on pixel $(i,j)$, 5×5 neighborhoods are established, as shown in Fig. 2. The neighborhood is divided into two regions according to different directions $\theta$. Four directional angles are used to ensure edge diversity.



Fig. 1. Flowchart of edge detection based on ACA.

Fig. 2.   Regional gray mean difference.

The mean value $B_{\theta_n}^{r_m}$ of regional grayscale can be expressed as Eq. (5).

$$B_{\theta_n}^{r_m} = \frac{\sum_{i,j \in D_m} f(i,j)}{N(i,j \in D_m)}, m = 1,2; n = 0,1,2,3 \tag{5}$$

In Eq. (5), $N(i,j \in D_m)$ represents the total pixels in the region $D_m$. $\theta_n$ stands for the direction angle. The mean difference $\Delta B$ of regional grayscale can be expressed as Eq. (6).

$$\Delta B = \max\left(\Delta B_{\theta_n}\right) \tag{6}$$

Based on the difference in regional grayscale mean $\Delta B$, the grayscale difference on both sides of the short line segment where pixel point $(i,j)$ is located is significant. It indicates that the position of the short line segment may be the boundary of the image. The 5×5 neighborhood can generate more short line segments compared with traditional structures, thereby avoiding rough estimation of line segments in the 7×7 neighborhood and improving the accuracy of boundary information. After removing the influence of noise points, a

fusion gradient $H(i,j)$ is constructed by combining the grayscale gradient with the average regional grayscale value, which has high edge detection accuracy. Eq. (7) displays the calculation method.

$$H(i,j) = x\psi_T(i,j) + y\Delta B(i,j) \tag{7}$$

In Eq. (7), $x$ and $y$ are all relationship coefficients, $x + y = 1$. To obtain the noise reduction effect and ensure the accuracy of edge computing, $x = 0.5$ and $y = 0.5$ are set as the weight of the difference between the gray gradient and the regional gray mean. By using stroke width transformation, the edge extracted image is transformed into a stroke width transform image containing stroke width information. Then it is integrated and classified based on stroke width information. The corresponding connected domains are extracted. On the left side of Fig. 3 is a typical stroke image. The red dots represent pixels and white dots represent the background. The boundary extraction of the original stroke map can obtain the stroke boundary map, as shown on the right side of Fig. 3.

### B. Text Extraction and Recognition Based on Binary Processing

After accurately locating the position of characters in the image using edge detection ACA, the determined character area is segmented and recognized. In the template library, there are only single characters. The existing methods can only recognize single characters. Therefore, to accurately recognize characters, the image needs to be segmented [22]. To grayscale a color image, the grayscale of each point in the image is changed to 0 or 255. The processed color image is presented as a gray image. To binarize the obtained gray image, a specific algorithm is used to obtain a corresponding threshold. Then, the grayscale of each point in the image is compared with this threshold. It is divided into object and background. After image binarization, the grayscale of each pixel in the image is only 0 and 255, excluding other grayscale sizes. In the binarization, all points in the image with grayscale above the threshold are targeted, and the grayscale is set to 255. On this basis, points with a grayscale lower than this threshold are used as backgrounds. Their grayscale is set to 0. The binary segmentation method is shown in Eq. (8).

$$g(x,y) = \begin{cases} 255, f(x,y) \geq T \\ 0, f(x,y) \leq T \end{cases} \tag{8}$$



Fig. 3.   Construction of edge point pairs.

In Eq. (8), $f(x,y)$ stands for the pixel value at the coordinate $(x,y)$. $g(x,y)$ is the pixel value at the $(x,y)$ coordinate in the image processed by the binarization method. A low-pass filter preprocesses the collected images to reduce or eliminate noise. The algorithm determines the optimal threshold, ensuring that the image can be effectively segmented into the target and background at the boundary of the optimal threshold. All points in the image with grayscale greater than the calculated threshold are set to 255, and points below the calculated threshold are set to 0. The obtained image only has two colors, black and white, which means dividing the image into target and non-target regions to achieve image binarization.

Binarization is a very important image processing technique. Selecting the appropriate threshold is an important step in image binarization. The Otsu algorithm uses a special discriminant function to determine the threshold size of binarization. The binarization algorithm divides the grayscale values of points on the entire image into two types at the threshold $t$. $C_0 = \{0,1,2,\cdots,t\}$ represents the background area [23]. If the total pixels in the character area after grayscale processing are $N$, and each pixel has the highest grayscale level $L$, the grayscale size of the entire image is within $[0,L-1]$. The pixel with a grayscale value of $i$ is $n_i$. The probability calculation method for $i$ is expressed as Eq. (9).

$$p_i = \frac{n_i}{N} \tag{9}$$

Assuming the threshold is $T$, the grayscale is $[0,T-1]$ for the $C_0$ region, and the grayscale is $[T,L-1]$ for the $C_1$ region, the probability of $C_0$ and $C_1$ occurring is shown in Eq. (10).

$$\begin{cases} W_0 = \sum_{i=0}^{T-1} p_i \\ W_1 = \sum_{i=T}^{L-1} p_i = 1 - W_0 \end{cases} \tag{10}$$

The grayscale mean of S and $C_1$ can be expressed as Eq. (11).

$$\begin{cases} u_0 = \sum_{i=0}^{T-1} i p_i \\ u_1 = \sum_{i=T}^{L-1} i p_i \end{cases} \tag{11}$$

The inter class variance $\sigma^2$ is displayed in Eq. (12).

$$\begin{cases} \sigma^2 = W_0 (u_0 - u)^2 + W_1 (u_1 - u)^2 \\ u = W_0 u_0 + W_1 u_1 \end{cases} \tag{12}$$

In Eq. (12), $u$ represents the average grayscale. $T$ increases in steps of 1 within the range of $[0,L-1]$. When $\sigma^2$ is the maximum, the corresponding $T$ value is the optimal binarization threshold.

The conventional Otsu algorithm often has unsatisfactory results when conducting binarization to images with slight differences between the target and background. Moreover, there are many parts that belong to the key regions of characters that have not been extracted. This study improves the effectiveness of the algorithm by adjusting the size of the threshold. The threshold size that needs to be adjusted for binarization is directly proportional to the average grayscale. The relationship between the required adjustment size $w$ and the average grayscale value $E$ is shown in Fig. 4.



Fig. 4. Relationship between adaptive fine-tuning amount $w$ and average gray value $E$.

In Fig. 4, the curve in the figure is an arc. Model 1 is a monotonically smooth curve fitting that is convex upwards. Model 2 is fitted with a monotonic smooth curve that is concave downwards. The downward concave curve fitting effect is better, that is, mathematical model 2. The following is a specific analysis for mathematical Model 2. Assuming $(x_0, y_0)$ is the center of a circle, the relationship between $x_0$ and $y_0$ can be expressed as Eq. (13).

$$y_0 = -\frac{E_{\max} - E_{\min}}{w_{\max} - w_{\min}} \left( x_0 - \frac{E_{\max} + E_{\min}}{2} \right) + \frac{w_{\max} + w_{\min}}{2} \tag{13}$$

In Eq. (13), $E_{\max}$ and $E_{\min}$ respectively represent the maximum and minimum of the average grayscale values for a single column in the image. $w_{\max}$ and $w_{\min}$ represent the maximum and minimum of the required adjustment amount in the entire image, respectively. $x_0$ determines the curvature of a circle. If $x_0$ is small, then the curvature is small. If $x_0$ is large, then the curvature is large. To ensure a moderate curvature of the fitted curve, the value between $E_{\max} - 50$ and $E_{\min} - 60$ is the most appropriate for $x_0$. The equation for the final fitted circle is shown in Eq. (14).

$$(w - x_0)^2 + (E - y_0)^2 = (x_0 - E_{\min})^2 + (y_0 - w_{\min})^2 \tag{14}$$

Therefore, the relationship between the fine-tuning amount $w$ and the average grayscale value $E$ of any column can be

expressed as Eq. (15).

$$w = -\sqrt{\left(x_0 - E_{\min}\right)^2 + \left(y_0 - w_{\min}\right)^2 - \left(E - y_0\right)^2 + x_0} \quad (15)$$

The relationship between fine-tuning amount and average grayscale value can be obtained from Eq. (15). On this basis, the traditional Otsu algorithm is improved using the above formula. The threshold after binarization is adjusted to achieve ideal segmentation results.

## IV. The Effectiveness Analysis of Text Recognition Methods

The experiment was carried out on an Intel i5-8250U 1.6$GHZ$ processor and 8GB of memory. The operating system was Windows 10. Matlab 2015b served as the operating platform. The algorithm performance was validated through the ICDAR2017RCTW database, including edge detection performance, text region localization performance, binarization performance, and image text recognition rate.

### A. Analysis of Character Positioning Effect

The research was based on the ICDAR2017RCTW database. Color text images and texture physical background images that meet the requirements were used to construct a data set. The text recognition effect was compared. The model was trained using TensorFlow 1.3. The programming language was Python. The server configuration was NVIDIA TESLA P4 graphics card. Fig. 5 was one of the images in the data set. It was used for subsequent recognition analysis.

To demonstrate the advantages of the Mallat wavelet fast decomposition algorithm in license plate edge detection, this study compared the improved algorithm with traditional edge extraction algorithms. Fig. 6 showed the comparison results. Fig. 6 (a) displayed the edge extraction effect of traditional algorithms. The text area was successfully separated from other areas. There was a clear difference between the white rectangular area and other connected areas. The text area was longer than other areas. Fig. 6 (b) showed the edge extraction effect of the improved algorithm. It could effectively preserve the edges of characters in the image. However, the traditional algorithms were prone to losing edge information of character regions during edge detection.

Fig. 7 showed the results of using the proposed algorithm to locate text regions in the image. From the figure, the improved algorithm located a total of seven main text regions in the image, basically including all the text in the image. The text regions could be well distinguished, providing a significant non-interference edge extraction effect.



Fig. 5. Original image.



(a) Edge detection graph



(b) Edge detection graph based on ant algorithm

Fig. 6. Comparison of processing effects of edge detection algorithms.



Fig. 7. Character positioning effect.

In order to test the effectiveness of ACA for text localization, the research selected four typical images from the data set for text detection and localization, and analyzed them through five indexes including correct detection rate, missed detection rate, false detection rate, recall rate, and precision rate. The positioning results are shown in Table I. Among them, the number of key frames and the total number of text lines were the total character area of the image. The number of correct detected lines and the number of missed lines represented the positioning accuracy of the algorithm, and the recall rate and precision rate represented the recognition accuracy of the algorithm. From Table I, the text localization effect of the proposed method was good. The text recall rate reached 95%, and the accuracy rate was around 85%. The accuracy rate was relatively low due to complex background interference, which couldn't exclude all factors, but it is already higher than traditional detection methods. The extracted text features were targeted. The ACA was effective, and the regional posterior method also reduced the false detection rate.

To verify the advantages of the proposed algorithm in processing text information detection, the study compared the edge point pair construction effects of the two algorithms. The results were shown in Fig. 8. From the projection of Fig. 8 (a) and Fig. 8 (b), the text area was more prominent because the irrelevant backgrounds were removed. It could effectively improve the character recognition rate in complex backgrounds and filter out interference in the background, thereby extracting characters more efficiently. In summary, this method could effectively solve the image segmentation in complex backgrounds, effectively reducing the impact of non character regions on characters, and reducing the character stroke missing.

### B. Analysis of Text Recognition Effect

The method of combining vertical and horizontal projection was used to segment the license plate area of vehicles. Taking Fig. 5 as an example, first, the license plate area was projected vertically. There were no peaks in the middle of each character. Based on this feature, the license plate was divided into individual characters to obtain a vertical text projection image, as shown in Fig. 9. From Fig. 9, the character area had more peaks, while the non-character area had no peaks. The area where the text was located had many peaks and valleys, and the peaks and valleys changed frequently. Therefore, the peak of the vertical projection curve had significant changes.

The study compared the binarization effects of text images. The binarized images obtained after processing with the Otsu algorithm were shown in Fig. 10. Fig. 10 (a) showed the Otsu binarization effect of intra class variance. In the figure, for color text images, the Otsu binarization effect of inter class variance was not very ideal. There were phenomena such as unclear characters in the text area. Fig. 10 (b) showed the Otsu binarization effect of intra class variance. Compared to the inter class variance algorithm, there was no significant improvement in the performance of color text images. Moreover, this algorithm was more complex than the inter class variance Otsu algorithm, and the computational efficiency significantly decreased. Fig. 10 (c) made appropriate adjustments to the threshold obtained by the Otsu algorithm. On the basis of the above threshold, the threshold was reduced by 15. In the figure, the binarization effect after reducing the threshold by 15 was better than before, especially for the binarization effect of the digital part.

TABLE I.    CHARACTER EXTRACTION RESULTS

| Index | Key frame count | Total lines of text | Correctly detect the number of rows | The missed rows | Number of misdetected rows | Recall factor | Precision ratio |
|---|---|---|---|---|---|---|---|
| Image 1 | 4 | 1 | 1 | 0 | 0 | 100.0% | 100.0% |
| Image 2 | 23 | 18 | 17 | 1 | 3 | 94.4% | 85.0% |
| Image 3 | 14 | 10 | 10 | 0 | 3 | 100.0% | 77.0% |
| Image 4 | 24 | 15 | 14 | 1 | 2 | 93.3% | 87.5% |
| Total | 65 | 44 | 42 | 2 | 8 | 95.3% | 84.0% |



(a) SWT algorithm projection diagram    (b) Improved algorithm projection diagram

Fig. 8.    Comparison of algorithm effects before and after improvement.

Fig. 9.  Projection of license plate in vertical direction.



(a) In-class variance Otsu binarization rendering



(b) Otsu binarized rendering of inter-class variance



(c) Threshold reduction by 15 binary effect

Fig. 10.  Comparison of binarization effects.

The images in the data set were located and the text area map was binarized, followed by text recognition. The recognition results were displayed in Table II. The recognition rate of Image 1 was 70.27%, which may be due to the low image clarity. The text recognition rate of the improved ACA proposed in the research generally reached over 80%, verifying the effectiveness of the binarization method. There was still room for further improvement in recognition rate. Overall, it had certain advantages compared to other algorithms.

TABLE II.     TEXT RECOGNITION RESULTS

| Index | Total number of characters in the location area | Correct word count | Recognition rate |
|---|---|---|---|
| Image 1 | 37 | 26 | 70.27% |
| Image 2 | 356 | 293 | 82.3% |
| Image 3 | 261 | 228 | 87.36% |
| Image 4 | 348 | 291 | 83.62% |
| Total | 1002 | 836 | 83.63% |

## V.  RESULTS AND DISCUSSION

Shot detection is the premise of key frame extraction (i.e., subtitle frame), and subtitle frame detection is the basis of subsequent subtitle positioning and recognition [24]. The wavelet decomposition algorithm is used to divide the video image into high frequency part and low frequency part. The high frequency part is the region where the characters are located. By comparing the wavelet decomposition algorithm with the traditional Sobel operator edge detection algorithm, it can be seen that the text region obtained by the wavelet fast decomposition algorithm is more complete. The localization of text area is divided into two processes, namely preliminary localization and precise localization. The character part is extracted from the whole video image by using the open and close operation function of MATLAB toolbox. The precise positioning of the character part refers to removing the edge of the extracted image character region, leaving only the part containing characters. With horizontal and vertical projection methods, using the character area projection gray histogram crest and trough continuous jump characteristics, the character is accurately located. As described in the review, the proposed algorithm can accurately determine the region where the characters are located in the video image, remove other background regions in the image, and the processed character region is clearly visible, and the noise is significantly reduced, which is conducive to the next step of character research.

After the text of the target area is located, in order to extract effective text information, the semantic part of the text needs to be extracted from the complex background, and the threshold segmentation of the text field is carried out again [25]. Subtitle extraction is actually a binary process, which requires that the resulting graph can basically maintain the original character features, without blank, hollow, and broken. Character extraction is mainly divided into two parts, including picture binarization and character segmentation. In

this paper, an improved Otsu algorithm is proposed for the binarization of character images, and the threshold is fine-tuned on the basis of the traditional Otsu algorithm. The amount of fine tuning is determined by the size of the average gray value. There is a curve function relationship between the average gray value and the fine-tuning amount [26]. Comparing the improved binarization effect with the traditional binarization effect, the improved binarization effect is significantly improved. The character area is clearer, the contrast with the background area is significantly improved, and the improved Otsu algorithm has higher processing accuracy under different lighting conditions. In the character segmentation stage, the research combines vertical projection and horizontal projection, and uses the characteristics of continuous peak value of character area and no white pixels in non-character area to remove the background area around characters in the image, so as to accurately locate the region where characters are located. In the character recognition stage, a character recognition algorithm based on template matching is used. The results showed that the text recognition rate of the improved ACA generally reached more than 80%, which verifies that the binarization method in this paper is effective.

## VI. CONCLUSION

ACA has the advantages of random search, fast convergence speed, and strong local search ability. This algorithm is used to binarize text images. The study combines the pheromone update mechanism and path selection mechanism in traditional ACA. The sensitivity to initial pheromone concentration and path selection strategies is improved. This to some extent enhances the local and global search capabilities. Then, the threshold obtained from the traditional Otsu algorithm is adaptively fine-tuned. The functional relationship between the fine tuning amount and the average grayscale value is found, completing the binarization of the image. The improved ACA could effectively preserve the edges of characters in images. Traditional algorithms were prone to losing 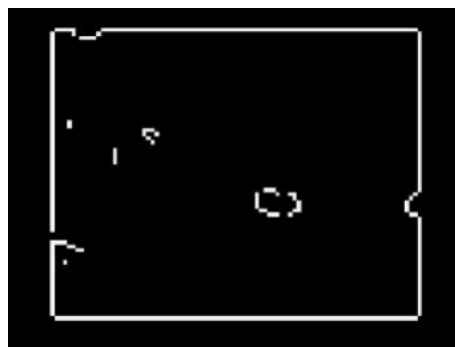edge information of character regions during edge detection. The text localization effect based on the improved ACA was good, with a text recall rate of 95% and a precision rate of about 85%. The regional posterior method also reduced the error detection rate. The binarization effect after reducing the threshold by 15 was better than before, indicating that fine-tuning the threshold improved the binarization effect. The text recognition rate of the improved ACA proposed in the research reached over 80%, verifying the effectiveness of the binarization method. Although the ACA is used in this study to avoid the positive and negative sample selection training of the classifier, the iterative nature of the ACA is bound to reduce the efficiency of the algorithm. In the future, it is necessary to study how to improve the efficiency of the algorithm, or develop an alternative method.

## ACKNOWLEDGMENT

## REFERENCES

[1] Cao H, Wu Y, Bao Y, Feng X, Wan S, Qian C. UTrans-Net: A Model for Short-Term Precipitation Prediction. Artificial Intelligence and Applications. 2023, 1(2): 106-113.

[2] Wu Z, Zhao Y, Zhang N. A Literature Survey of Green and Low-Carbon Economics Using Natural Experiment Approaches in Top Field Journal. Green and Low-Carbon Economy, 2023, 1(1): 2-14.

[3] Marchi S D, Erb W, Marchetti F, Perracchione E, Rossini M. Shape-Driven Interpolation with Discontinuous Kernels: Error Analysis, Edge Extraction, and Applications in Magnetic Particle Imaging. SIAM Journal on Scientific Computing, 2020, 42(2):472-491.

[4] Xiao H, Xiao S, Ma G, Li C L. Image Sobel edge extraction algorithm accelerated by OpenCL. Journal of supercomputing, 2022, 78(14):16236-16265.

[5] Lin F, Yu Z, Jin Q. Semantic Segmentation and Scale Recognition–Based Water-Level Monitoring Algorithm. Journal of Coastal Research, 2020, 105: 185-189.

[6] Zhang J, Li C, Rahaman M M, Yao, Y D, Ma P L, Zhang J H, Zhao X, Jiang T, Grzegorzek M. A comprehensive review of image analysis methods for microorganism counting: from classical image processing to deep learning approaches. Artificial intelligence review, 2022, 55(4):2875-2944.

[7] Usman A M, Abdullah M K. An Assessment of Building Energy Consumption Characteristics Using Analytical Energy and Carbon Footprint Assessment Model. Green and Low-Carbon Economy, 2023, 1(1): 28-40.

[8] Gan J, Wang W, Lu K. Compressing the CNN architecture for in-air handwritten Chinese character recognition - ScienceDirect. Pattern Recognition Letters, 2020, 129(6): 190-197.

[9] Yi N, Xu J, Yan L, Huang L. Task optimization and scheduling of distributed cyber–physical system based on improved ant colony algorithm. Future Generation Computer Systems, 2020, 109(23): 134-148.

[10] Tang Z, Ma H. Path guidance method for unmanned vehicle based on improved potential field ant colony algorithm. International Journal of Vehicle Design, 2022, 89(2): 84-97.

[11] Zhu S, Zhu W, Zhang X, Cao T. Path planning of lunar robot based on dynamic adaptive ant colony algorithm and obstacle avoidance. International Journal of Advanced Robotic Systems, 2020, 17(3): 4149-4171.

[12] Yu J, You X, Liu S. Dynamically induced clustering ant colony algorithm based on a coevolutionary chain. Knowledge-based systems, 2022, 251(5): 336-351.

[13] Hwang W, Kang D, Kim D. Brain lateralisation feature extraction and ant colony optimisation-bidirectional LSTM network model for emotion recognition. IET Signal Process, 2021, 16(1): 45-61.

[14] Guptha N S, Balamurugan V, Megharaj G, Sattar K N A, Rose J D. Cross lingual handwritten character recognition using long short term memory network with aid of elephant herding optimization algorithm. Pattern recognition letters, 2022, 159(7): 16-22.

[15] Liu Z, Pan X, Peng Y. Character Recognition Algorithm Based on Fusion Probability Model and Deep Learning. The Computer Journal, 2020, 64(11): 1705-1714.

[16] Alam M S, Kwon K C, Kim N. Implementation of a Character Recognition System Based on Finger-Joint Tracking Using a Depth Camera. IEEE transactions on human-machine systems, 2021, 51(3): 229-241.

[17] Chandio A A, Leghari M, Kakepoto I, Leghari M, Hussain G. Multilingual Natural Scene Character Classification and Recognition with Histogram of Oriented Gradients. Fourrages, 2020, 243(9): 86-94.

[18] Lee S H, Yu W F, Yang C S. ILBPSDNet: Based on improved local binary pattern shallow deep convolutional neural network for character recognition. IET image processing, 2022, 16(3): 669-680.

[19] Liu Y, Gao T, Song B, Huang C. Personalized Recommender System for Children's Book Recommendation with A Realtime Interactive Robot. Journal of Data Science and Intelligent Systems, 2023, 2(1): 1-6.

[20] Liu Y, Ma X, Li X, Zhang C. Two-stage image smoothing based on edge-patch histogram equalisation and patch decomposition. IET Image Processing, 2020, 14(6):1132-1140.

[21] Kim B, Kim J, Ye J C. Task-Agnostic Vision Transformer for Distributed Learning of Image Processing. IEEE Transactions on Image Processing, 2023, 32(20):2023-218.

[22] Garai S, Paul R K, Kumar M, et al. Intra-Annual National Statistical Accounts Based on Machine Learning Algorithm. Journal of Data Science and Intelligent Systems, 2023, 2(2): 12-15.

[23] Huang B, Lin J, Liu J, Chen J, Zhang J, Hu Y, Chen E, Yan J. Separating Chinese Character from Noisy Background Using GAN. Wireless Communications and Mobile Computing, 2021, 54(1): 1-13.

[24] Ding C, Zhou A, Liu X, Ma X, Wang S G. Resource-aware Feature Extraction in Mobile Edge Computing. IEEE Transactions on Mobile Computing, 2020, 21(1):321-331.

[25] Clausner C, Pletschacher S, Antonacopoulos A. Flexible character accuracy measure for reading-order-independent evaluation - ScienceDirect. Pattern Recognition Letters, 2020, 131(5): 390-397.

[26] Wang R, Cao W, Wu S, Jia M, Wang X P. Optical character correction of large-curvature annular sector text in polar coordinate system. Pattern recognition letters, 2023, 167(5): 157-163.

# Application of Style Transfer Algorithm in Artistic Design Expression of Terrain Environment

Yangfei Chen

Department of Architecture, Sichuan College of Architectural Technology, Deyang, 618000, China

*Abstract*—The use of artistic expression to depict and express terrain and landform can not only convey terrain information, but also spread art and culture. The existing landscape design methods focus on the accurate expression of terrain height and the realistic expression of form, but neglect the aesthetic aspect of landscape design. In view of this situation, this paper studied the use of generative adversarial network, constructed the presentation mode of landscape plane style, and realized the expression of landscape art style. A terrain style transfer model based on a pre-trained deep neural network model and style transfer algorithm was constructed to achieve a variety of terrain style expressions. The results showed that, in terms of Peak Signal to Noise Ratio, the proposed style transfer algorithm was higher than style attention network and adaptive instance normalization, and the peak signal-to-noise ratio index value was increased by 7.5% and 16.5%. This indicated that the style transfer model proposed by the topographic artistry research had more advantages in terms of image diversity and fidelity. The Structural Similarity Index of the proposed algorithm has been greatly improved. This research expands the method of computer rendering of terrain environment art, which is of great significance for the preservation of traditional Chinese culture.

*Keywords*—*Generative adversarial network; terrain; style transfer; artistic; peak signal-to-noise ratio; Structural Similarity Index*

## I. INTRODUCTION

Geographic abstraction is one of the important methods of geographical cognition. It is a key feature obtained by summarizing the geographical reality. There are two kinds of abstract processes on the basis of scientific cognition of topography: scientific abstraction and artistic abstraction [1-2]. Each of these topographic abstract processes has its own advantages and characteristics. Each formed a perfect and mature theoretical system, which can serve people's cognition of topographic landform. On the basis of terrain abstraction, the presentation of terrain expression is a problem worth studying. For a long time, domestic and foreign scholars have been committed to the quantitative and fine expression of terrain. The research on scientific and accurate expression of terrain has developed rapidly [3-4]. From hachure method, hill shading method, contour line method and other terrain three-dimensional scene construction, people can use computer to simulate the realistic terrain and ground objects. The scientific expression of terrain abstraction is becoming more and more mature. However, science and aesthetics are a pair of contradictory and unified contradictions. The aesthetic examination in science and the scientific connotation in aesthetics promote each other [5-6]. The existing research is more focused on the scientific and accurate expression of terrain, which cannot fully meet the diverse visual requirements of people. The study analyzes the representation of terrain through artistic expression by incorporating painting techniques to enhance the aesthetic quality of terrain representation. This research has significant implications for the study of artistic terrain representation. By using computer related technology, the expression technique of painting art with national characteristics is combined with the expression method of terrain to realize the artistic expression of terrain style. While enriching the terrain to express the cultural connotation, it can also make the reader feel the artistic charm of Chinese painting. Based on the in-depth study of terrain expression and artistic style learning, this paper first determines the extraction method suitable for terrain expression and forms the dataset of style samples. Then, a style transfer model oriented to terrain expression is constructed. Finally, the real terrain structural feature information is taken as the model transfer target to achieve the artistic expression of terrain style. This approach effectively integrates scientific and artistic expression of terrain expression. The research results are expected to enrich the content of terrain expression and improve the effect of artistic expression.

This study is mainly separated into the following six sections. Section II is a literature review on intelligent algorithms and style transfer applications. Section III is the construction of a style transfer model on the ground of generative adversarial networks (GAN). Section IV is an analysis of the algorithm performance and application results of the style transfer model proposed by the research institute. Results and discussion is given in Section V. Finally, Section VI concludes the paper.

## II. RELATED WORK

The research on traditional graphic art style can be roughly summarized into three categories: texture synthesis based on sample images, virtual brush model building based on strokes, and local separation processing based on image structure. Zhang and other scholars used GANs combined with attention mechanisms to convert real-world facial and cartoon style images into unpaired datasets. The results show that the proposed new model network avoids the complexity of the model and achieves a good balance in the task of style and content transformation [7]. Wang et al. proposed a cyclic consistent GAN on the ground of edge features and self attention for style transformation in visual effects. The model structure includes a generator, a discriminator, and an edge feature extraction network. The results show that the model has advantages in style conversion, as it can better preserve

the details of the original image and has good image quality [8]. To promote the research on image style transfer on the ground of neural networks, Huang and other research teams summarized and discussed the main principles and methods of image style transfer on the ground of neural networks, and provided a detailed description of neural networks on the ground of neural networks [9]. Scholars such as Zhou proposed an image based convolutional neural network transfer model on the ground of deep mixing. This deep hybrid generation model mainly relies on the combination of adversarial network generation and self encoder [10]. Wei et al. proposed a transformation based visual style transfer method for rendering style patterns. The results indicate that under appropriate style supervision, the transformer can learn similar texture bias features like CNN [11]. Gupta and other research teams used transfer learning to pre train convolutional neural networks for image style transfer tasks. Using these models can generate high perceptual quality images, which are a combination of the content of any image and the appearance of famous artworks [12]. Kim and other scholars proposed a CNN inference accelerator for style transformation applications, which utilized network compression and layer chain technology. In addition, layer chain technology has been proposed to reduce off chip memory traffic, thereby increasing throughput at the cost of smaller hardware resources. The results show that this method has good performance in changing the style of content images [13].

In the field of environmental design, Yang and other scholars reviewed the development process of environmental art and design discipline, relevant key educational institutions, and the opinions of well-known scholars. In addition, in the context of emerging design trends, new collective and constructive environmental art design ideas and methods have also been proposed. Finally, this study suggests the uniqueness and regionality of design culture in the Chinese context through the study of the name dispute for environmental art design in China [14]. Scholar Wang conducted research on computer-assisted interaction of vision, elaborating on digital technology and its development trends in the future, and also pointing out that technology is the future practice of digital media art. It proposes a new form of art, which is the "intelligent visual art" that conforms to the development of the intelligent era, and explores the opportunities and challenges of intelligent visual art [15]. Liu and his research team conducted two user studies to evaluate the impact of three layouts with different curvatures around users in a virtual environment (flat wall, semi circular surround, and circular surround) on visual spatial memory tasks. The results show that compared to the circular layout, participants were able to recall spatial patterns more accurately and report more positive subjective ratings [16]. Liu and other scholars proposed a multi-dimensional urban landscape design method on the grounds of nonlinear theory to solve the problem of large differences in multi-dimensional urban landscape design. On the grounds of the parameterized model method, multi-dimensional nonlinear landscape design has been implemented, improving the quantitative analysis ability of multi-dimensional nonlinear landscape design [17]. Wang et al. proposed a virtual environment on the ground of virtual reality technology and intelligent algorithms, which uses a 3-bit binary to represent digital factors and creates a virtual environment by simulating the display environment. The research shows that applying virtual reality technology and intelligent algorithms to landscape design in coastal areas is feasible and has achieved certain results [18]. The purpose of research by scholars such as Thamrin is to achieve community service through collaborative design in interior design teaching. The article describes the learning and design methods on the ground of the human centered design method in collaborative design, and analyzes the benefits brought by this method [19].

In summary, there are currently many studies using convolutional neural networks or GANs to construct style transfer models for text or images. In the research of environmental artistic design, more pure theories or suggestions are presented, and less intelligent algorithms are introduced to achieve the artistic expression of terrain. In view of this, this study constructs a style transfer model for terrain artistic expression on the ground of GANs and two novel channel models.

## III. RESEARCH ON STYLE TRANSFER ALGORITHM FOR TERRAIN REPRESENTATION

To address the issues of unsatisfactory performance of art attributes and incompatibility with style characteristics caused by image style conversion, this study aims to construct an image style conversion network that enhances art style attributes. On this basis, feature transfer and refinement of the two pathways are achieved through the attribute refinement pathway of the style domain and the semantic information refinement pathway of self attention. The former can highlight the aesthetic characteristics of the image, while the latter can combine the style characteristics of the image with its most suitable content characteristics on the grounds of the characteristics of the image. A pooling layer on the grounds of wavelet based multi-level decomposition has been introduced into the compiler of the generation algorithm, which maximizes the preservation of important features of the image during feature transfer. On this basis, a dual resolution discrimination network is used to distinguish different types of images, to obtain images similar to different types of images.

### A. Extraction of Terrain Environmental Features

In the dissemination of style information, using a single channel style transformation branch is difficult to retain its unique style domain characteristics. This is mainly because after encoding the image content and style, there is no mapping between the semantic and style styles. However, using direct fusion during the decoding process will bring significant errors, which will prevent the decoding system from efficiently training attributes such as color and texture. This can reduce the decoding system's ability to extract deep style features, ultimately resulting in changes in the strokes in the generated graphics, i.e. missing attributes in the style domain. In response to the requirement to focus on style domain attributes in style conversion, this study proposes to use region perceptrons as the extraction channel for style domain attributes. On the grounds of a comprehensive analysis of texture characteristics, high-precision extraction in the style domain is achieved and converted into style domain

attributes. Fig. 1 is a schematic diagram of the style domain attribute refinement channel.



Fig. 1.    Style domain attribute refinement channel.

It obtains image feature maps at different levels through encoders at different levels, and obtains texture and color information of the image through the Gram matrix. In the process of encoding the style, the feature maps of different levels are segmented into multiple channels through 3x3 convolution, and then wavelet fusion operation is added to obtain the feature maps of the channel set, thereby avoiding the loss of local parameters without increasing the size of the feature maps. Finally, a style feature map on the ground of channel sets is constructed, and a weighted convolution method is used for forward diffusion to obtain the style domain attributes of each level, thereby achieving complete migration and accurate fusion. Eq. (1) is the calculation expression for the feature map.

$$F_i^{DI} = h\left(\left[FC_i\left(Gram\left(F_i\right)\right) \otimes F_i'\right]\right) \tag{1}$$

In Eq. (1), $F_i$ represents the feature maps obtained from different layers. $F_i'$ represents the channel set feature map. $\alpha_i$ represents the style domain attribute. $\otimes$ represents the channel connection operator. $FC$ is the fully connected layer. $h$ is the weight sharing fully connected layer. Eq. (2) calculates an expression for the style domain attribute.

$$\alpha_i = \sigma\left(F_i^{DI}\left(I\right)\right) \tag{2}$$

In Eq. (2), $\sigma$ represents the sigmoid operation. By refining the attributes of the style domain between the modules that compile the style image, a communication bridge is established to match similar types of features, thus achieving a direct transfer of style domain attributes. This reduces the error of style feature transfer and solves the

problem of generating style images without style domain attributes. Meanwhile, through this new approach, the instability of overall feature transfer is enhanced, and the overall style information such as color and texture is effectively integrated to achieve overall style transfer, thereby achieving a more harmonious visual effect. Due to the significant regional distribution characteristics of patterns and the differences in style domain information displayed by different regions, it is necessary to achieve synchronous maintenance of the overall and local styles, and to adaptively adjust the constructed neural network to achieve automatic matching of the most consistent style features. This method utilizes a self attention mechanism to train the regularized model, and improves it to have higher credibility and average accuracy. This is to achieve accurate matching of images with the same semantics in type images. This is precisely in line with the process of image migration, which involves matching the content contained in the image with the semantic information contained in the image to achieve the closest possible style transfer. Fig. 2 shows the self attention semantic feature matching channel.

Using factor decomposition paradigm standardization, standardized text information and style information are obtained as input through self focused semantic feature matching channels. On this basis, the study extracted the semantic consistency between the two and their key features, thus achieving the learning of standardized parameters. This method can achieve dynamic pixel by pixel shift and scaling of content features, while maintaining content features to match the style characteristics of the text. The conversion process can be expressed by Eq. (3).

$$\overline{F_c'}, \overline{F_s'} = \gamma_s \otimes \left(\alpha_{c\_s}\left(\overline{F_c}, \overline{F_s}\right) + \alpha_{c\_s}\right) + \beta_s \tag{3}$$

Fig. 2. Self attention semantic feature matching channel.

In Eq. (3), $\gamma_s$ and $\beta_s$ represent normalization parameters respectively. $\overline{F_c}$ and $\overline{F_s}$ represent normalized content and style features, respectively. $\overline{F_c'}, \overline{F_s'}$ represent the content and style features refined by the self attention semantic feature matching channel. $\alpha_{c\_s}$ and $\alpha_{s\_c}$ are unit vectors along the channel dimension, which obtain standardized content and transition information between styles before style matching. Eq. (4) represents the normalization parameter calculation formula.

$$\begin{cases} \gamma_s = ReLU\left(conv_{1\times1} \otimes SA_\gamma\left(\overline{F_c}, \overline{F_s}\right)\right) \\ \beta_s = ReLU\left(conv_{1\times1} \otimes SA_\beta\left(\overline{F_c}, \overline{F_s}\right)\right) \end{cases} \quad (4)$$

$SA_\gamma$ and $SA_\beta$ represent self attention convolutional network layers.

*B. Construction of Style Transfer Model*

On the grounds of the GAN, a picture style transmission mode was established to enhance artistic style attributes. The function of the generator is to take content images and style images as inputs, responsible for representing mapping relationships and obtaining the generated images [20]. Discriminator is a dual resolution discriminator network that can distinguish between original and newly created images by using the original image as the correct sample and the generated image as the error sample. Fig. 3 is a schematic diagram of the generator structure.

This study proposes a generator network consisting of an encoder, a dual channel characteristic transmission module, and a decoder. The dual channel feature transmission module consists of two parts: one is the type field attribute refinement path. The other part is the self attention path. On the grounds of content images and style images as inputs, novel style images are created through the structural features of content images and the style features of style images. In addition, to maintain the edge features and style features of the content image, a new pooling layer has been introduced in the codec, as shown in Fig. 4.

The wavelet multi-layer transformation network decomposes feature maps at multiple levels, achieving effective extraction of feature maps in high and low frequency bands. In the image, the parts with high color impact exhibit a stepped change in grayscale, which best displays the details of the image and is also a characteristic of the image. On this basis, the high-frequency sub band information in the region is matched with it, and then synthesized into a feature band, which is then converted into a sub feature band by the excitation function of the band. For the background area of the image, its grayscale distribution is within a certain range, and the differences in values are not significant, occupying the largest part of the entire image area. It represents the low-frequency subbands of the image, which can be used as the input for each level of filter. Eq. (5) is the mathematical expression for this process.

$$\{LL, LH, HL, HH\} = F_{DWT}(LL) \quad (5)$$

In Eq. (5), $(HH, LH, HL)$ represents high-frequency subband information. $LL$ represents low-frequency subband information. $f_1$ represents the feature map. $F_{DWT}$ represents the filter operation of wavelet transform. Eq. (6) is the calculation ratio expression for feature map $f_1$.

$$f_1 = ReLU\left(BN\left(Conv_{3\times3}\left(LH, HL, HH\right)\right)\right) \quad (6)$$

In the reconstruction step, it reverses the decomposition step by upsampling the sub feature maps, then convolves them with a filter to obtain the sub feature maps, and finally obtains the reconstructed feature maps. Finally, this information is input into the encoder/decoder for the next step of feature extraction. Fig. 5 shows the discriminator network structure.

Fig. 3. Schematic diagram of generator structure.



Fig. 4. Pooling layer based on wavelet multilevel transform.



Fig. 5. Discriminator network structure.

The discriminator adopts a dual resolution type, which distinguishes the generated images by combining high-definition discriminators with traditional ones. The high-definition discriminator first enlarges the image to 512x512, and then inputs the configuration of the discriminator into the discriminator. Improved resolution allows for better capturing of textures in images. Its main goal is to obtain more realistic images by imposing more restrictions on the high-resolution characteristics of the images. Traditional resolution discriminators can effectively constrain the global structure of the generated image and ensure the semantic correlation between images, thereby improving the quality of the image. Eq. (7) is the loss function of GAN network, style domain attribute transformation network, and feature transformation module.

$$L_{AAH-GAN} = E_x\left[InD(I_s)\right] + E_z\left[In\left(1 - D\left(G(I_c)\right)\right)\right] + L_{DI} + L_{SAFIN} \quad (7)$$

In Eq. (7), $G(I_c)$ represents the generated stylized image. $D(I_s)$ represents the probability that the input image will be judged as true. $D(G(I_c))$ represents the probability that the stylized image given by the discriminator is a real image. $L_{DI}$ represents the loss of style domain attribute refinement channels. $L_{SAFIN}$ represents the loss of self attention semantic feature matching channel. Eq. (8) refines the channel loss function expression for style domain attributes.

$$L_{DI} = \lambda_{bce}L_{bce} + \lambda_{dlow}L_{dlow} \quad (8)$$

In Eq. (8), $L_{bce}$ represents domain attribute loss. $L_{dlow}$ represents the loss of feature fusion process. $\lambda_{bce}$, $\lambda_{dlow}$ represent the weighting factor for domain attribute loss and feature fusion process loss. Eq. (9) is the fusion feature calculation formula.

$$I_{mix} = mix\left(F_{s_c}, F_{s_s}, \beta\right) \quad (9)$$

In Eq. (9), $I_{mix}$ represents the fused features. *mix* represents the Mixup method. $\beta$ represents the interpolation strength. Eq. (10) is the calculation formula for the channel loss function of self attention semantic feature matching.

$$L_{SAFIN} = L_c + \lambda_s * L_s \quad (10)$$

In Eq. (10), $L_c$ and $L_s$ represent loss of content and style. $\lambda_s$ represents the weight coefficient assigned to style loss relative to content loss.

## IV. UTILITY ANALYSIS OF IMPROVING GAN'S TERRAIN ARTISTIC STYLE TRANSFER MODEL

To verify the effectiveness of the GAN style transfer model on the grounds of dual channel improvement proposed by the research institute, this study divided the experiment into performance verification and application effect verification stages. Performance testing mainly evaluates the algorithm itself, while the application effect is related to the style transfer of terrain images.

### A. Performance Analysis of Style Transfer Model on the Ground of Dual Channel Improved GAN

The system used in the experiment was UBUNTU-18_ cuda10.1, which was accelerated using two Ge ForceRTX2080 Ti chips. On this basis, an Adam optimizer was used with a learning rate of 0.0001 and an iteration number of 2000 as initial parameters. The dataset contains 6000 Impressionist images, all of which can be resized to 256x256, and any size image can be used during the testing period. The comparative algorithms used in the study are CycleGAN, StyleGAN, and Pix2Pix. The evaluation indicators are accuracy, recall, and F1 value.



(a) Accuracy results of different models on the test set

(b) Accuracy results of different models on the training set

Fig. 6. Accuracy bar chart.

Fig. 6 (a) shows the training results of the model on the test set, the proposed model in the study achieved the highest accuracy values, corresponding to the accuracy results of 93.8%, 95.1%, and 94.6% for the three data groups, respectively. Fig. 6 (b) shows the training results of the model on the training set. The StyleGAN and Pix2Pix algorithms have the highest accuracy of 92.8% and 93.7% on the three data sets. This indicates that the highest accuracy of the model proposed in the study has been improved by 2.3% and 1.4%, respectively, compared to the latter two algorithms.

Fig. 7 (a) and Fig. 7(b) show that the proposed algorithm has the highest recall rate on both the test and training sets. On the test set, the recall rates of the models proposed by the research institute are all above 90%, with a maximum value of 94.6%, while the highest recall rates corresponding to StyleGAN and Pix2Pix algorithms are 90.2% and 87.4%, respectively. By comparison, it can be seen that the algorithm proposed by the research institute has improved the recall rate by 4.4% and 7.2%.



(a) Recall results of different models on the test set



(b) Recall results of different models on the training set

Fig. 7. Recall rate box plot.



(a) Training set



(b) Test set

Fig. 8. F1 value line chart.

Fig. 8 shows that the model proposed in the study is at its highest position on the F1 value curve. Fig. 8 (a) shows the training results of the model on the test set, when the number of iterations is 100, the F1 value converges to 97.6%, which is 2.3% higher than the StyleGAN algorithm. Fig. 8 (b) shows the training results of the model on the training set, when the iteration is completed, the F1 value of the proposed algorithm converges to 95.3%, which is higher than 95%.

### B. Analysis of the Application Results of Style Transfer Model Based on Dual Channel Improved GAN

The experimental dataset remains unchanged, and the

comparative algorithms used in the study are Adaptive Instance Normalization (AdaIN), Style Attention Network (SANet), CycleGAN, and SAFIN. Using Peak Signal to Noise Ratio (PSNR) to measure the distortion state, the size of PSNR reflects the similarity between the migrated image and the reference. The Structural Similarity Index (SSIM) parameter is used to measure the similarity between the content image and the generated image. The closer SSIM is to 1, the more similar the structure is. Using IS scores to evaluate the clarity and diversity of generated images, the higher the IS score, the better the quality of the generated images. Fig. 9 shows the PSNR curves of different algorithms.

Fig. 9. PSNR curves of different algorithms.

Fig. 9 (a) shows that the PSNR curves of the images obtained from the training dataset exhibit consistency, all gradually increasing and tending to converge. As the iteration process progressed, the proposed algorithm achieved the highest PSNR index value of 95.9 when the number of iterations was 80. The SANet algorithm achieved a second highest PSNR index value of 88.7 when the number of iterations was 85. The AdaIN algorithm achieves a minimum PSNR metric value of 80 when the number of iterations is 90. By comparison, it can be seen that the algorithm proposed in the study achieved the fastest convergence speed and the highest PSNR index value, which were 7.2 and 15.9 higher than the latter, respectively. Fig. 9 (b) shows that the proposed algorithm achieved the highest PSNR index value in the test set, with a size of 71.8, and the overall PSNR curve remained around 70 with minimal fluctuations. The overall PSNR index value of the SANet algorithm fluctuates around 50. When the iteration reaches the middle and late stages, the curve fluctuation increases. When the iteration number is 100, the PSNR index value is 52.6, which is 19.2 lower than the model proposed in the study. The PSNR value curve of AdaIN algorithm is at its lowest point throughout the entire iteration time, and there is a certain downward trend in the early stages of the iteration. The final PSNR value is 40.

Fig. 10 (a) shows that on the training set, the algorithm proposed in the study is closest to 1 in terms of SSIM index. At the end of the iteration, the corresponding SSIM index value is 92.5%. And the curve begins to converge when the number of iterations is 30, with almost no fluctuations in the convergence process and excellent convergence performance. The SSIM index curves corresponding to the SANet algorithm and AdaIN algorithm both have two climbing processes, with 64 and 80 iterations starting to converge, both of which are higher than the algorithm proposed in the study. In terms of SSIM convergence value, the SSIM values of the latter two algorithms are 89.9% and 82.6%, which are reduced by 2.6% and 9.9% compared to the algorithm proposed in the study. Fig. 10 (b) shows that on the test set, the SSIM index curve of the proposed algorithm in the study shows a continuous climbing trend, with a relatively smooth upward process and no falling curve segments. When the number of iterations is 100, the SSIM index value of this curve is 93.6%, which is higher than 90.0%. The SSIM curves of the SANet algorithm and AdaIN algorithm both show a downward trend during the iteration process, and the overall SSIM index values are both below 80%. At the completion of the iteration, the SSIM index values of the two algorithms were 72.8% and 55.6%, respectively. By comparison, it can be seen that the SSIM index values of the algorithm proposed in the study have increased by 20.8% and 38% compared to the latter two algorithms.



Fig. 10. Different algorithm SSIM curves.

Fig. 11. Different algorithm IS curves.

Fig. 11 (a) and Fig. 11 (b) show that the IS score curves of the three algorithms show the same trend of change, with a slow upward trend in fluctuations, regardless of whether it is the training set or the test set. Among them, the IS score of the algorithm proposed in the study has always been at the highest level during the iteration. When the number of iterations is 100, the IS scores on both datasets are 91.8% and 90.0%, while the AdaIN and SANet algorithms have IS scores below 90% on the training set and test set. By comparison, it can be seen that the algorithm proposed in the study has a higher IS score, indicating that the algorithm proposed in the study has a more naive clarity and diversity in image transfer.

## V. RESULTS AND DISCUSSION

The expression of terrain art features is subjective and complex, making it challenging to use existing art style learning algorithms for style rendering. In this paper, a deep convolutional neural network for distinguishing and extracting stylistic texture features is designed based on a pre-trained deep neural network model. A sample feature dataset containing different styles is made, which is used to train the network. Combined with the idea of style transfer algorithm, a terrain style transfer model is constructed based on deep neural network to realize the expression of multiple styles of terrain. Based on the observation of various scenes, different styles of terrain expression can be perceived. The experimental results show that the proposed method has the following characteristics: 1) The terrain style transfer model can generate different styles of terrain representation by adjusting content reconstruction and style reconstruction factors. 2) The model is suitable for the style rendering of a large range of terrain scenes. 3) It can realize the artistic expression of terrain ink painting style from multiple perspectives.

Terrain style expression is a form of terrain artistic expression that involves scientific cognition, artistic abstraction and artistic expression. It requires a balance between science and art. The results of the topographic expression of style were evaluated from the aspects of science and aesthetics. On the one hand, the validity of the result of terrain artistic expression is evaluated qualitatively. The results show that the topographic style transfer results generated by the research extraction method can take into

account the scientific expression (the maintenance of the topographic feature structure) and the artistic expression (the aesthetic expression of the style). The number of topographic feature elements affects the overall effect of topographic style transfer. The extracted topographic feature elements are too few or too many, and cannot show a good effect. Through the evaluation of multiple style renderings, it is shown that the result of the combination of multiple convolutional layers in a deep neural network produces better results in expressing the style image, resulting in a smoother and more continuous visual experience. On the other hand, using texture analysis method for reference, the transfer effect of style texture is quantitatively analyzed. It is found that the general trend of each texture feature parameter curve is similar to the style sample, which indicates that the model can accurately realize the transfer of multiple styles. The migration degree of style texture in the model is different in four aspects: the ranking from high to bottom is manifested as obvious degree > complexity degree > similarity degree > thickness degree.

## VI. CONCLUSION

Combining the expressive techniques of graphic art with the expressive techniques of landforms can enhance the beauty of landforms, thereby enriching their artistic expression. To achieve the artistic expression of the terrain landscape, this study was on the grounds of generating adversarial networks and establishing a new style transfer mode. This pattern included two channels: attribute refinement in the style domain and self-focused semantic feature matching. It utilized the pooling layer of wavelet multi-level transformation to preserve the boundary and style characteristics of the content image. The results showed that in terms of SSIM indicators, the corresponding value of the model was closer to 1, and the convergence value of SSIM was 0.925. The SSIM convergence values of the other two algorithms were 0.899 and 0.826, both of which are below 0.9. The algorithm proposed by the research institute exhibited the same iterative trend in IS score indicators as AdaIN and SANet algorithms. However, in the process of change, the model proposed by the research had a higher IS score, with an IS score of 91.8% at the completion of the iteration. This indicated that the terrain artistic style transfer model proposed by the research institute had more advantages in presenting image clarity and richness.

Therefore, this model had certain application potential in reconstructing style content for artistic expression of terrain. The scientificity was an important factor to evaluate the result of the style expression of terrain ink painting. However, this paper only provided a qualitative measurement. In the future, on the basis of this study, it is necessary to deeply analyze the results of topographic ink style transfer, and put forward a quantitative analysis method of topographic information preservation.

REFERENCES

[1] Klein J T. Boundary Discourse of Crossdisciplinary and Cross-Sector Research: Refiguring the Landscape of Science. Minerva, 2023, 61(1): 31-52.

[2] Szkola J, Piza E L, Drawve G. Risk terrain modeling: Seasonality and predictive validity. Justice Quarterly, 2021, 38(2): 322-343.

[3] Hocini N, Payrastre O, Bourgin F. Performance of automated methods for flash flood inundation mapping: a comparison of a digital terrain model (DTM) filling and two hydrodynamic methods. Hydrology and Earth System Sciences, 2021, 25(6): 2979-2995.

[4] Fang B, Jiang M, Shen J. Deep Generative Inpainting with Comparative Sample Augmentation. Journal of Computational and Cognitive Engineering, 2022, 1(4): 174-180.

[5] Lei Y. Research on microvideo character perception and recognition based on target detection technology. Journal of Computational and Cognitive Engineering, 2022, 1(2): 83-87.

[6] Choudhuri S, Adeniye S, Sen A. Distribution Alignment Using Complement Entropy Objective and Adaptive Consensus-Based Label Refinement for Partial Domain Adaptation//Artificial Intelligence and Applications. 2023, 1(1): 43-51.

[7] Zhang T, Yu L, Tian S. CAMGAN: Combining attention mechanism generative adversarial networks for cartoon face style transfer. Journal of Intelligent & Fuzzy Systems, 2022, 42(3): 1803-1811.

[8] Wang L, Wang L, Chen S. ESA-CycleGAN: Edge feature and self-attention based cycle-consistent generative adversarial network for style transfer. IET Image Processing, 2022, 16(1): 176-190.

[9] Huang K, Lei M, Zhou B, Yang Q. Research Progress on Deep Learning based Image Style Migration Methods. International Core Journal of Engineering, 2022, 8(11): 46-54.

[10] Zhoua J, Wangb Y, Gongc J, Dong G Q，Ma W T. Research on Image Style Convolution Neural Network Migration Based on Deep Hybrid Generation Model. Academic Journal of Computing & Information Science, 2021, 4(8): 83-89.

[11] Wei H P, Deng Y Y, Tang F, Pan X J, Dong W M. A comparative study of CNN-and transformer-based visual style transfer. Journal of Computer Science and Technology, 2022, 37(3): 601-614.

[12] Gupta V, Sadana R, Moudgil S. Image style transfer using convolutional neural networks based on transfer learning. International journal of computational systems engineering, 2019, 5(1):53-60.

[13] Kim S, Jang B, Lee J, Bae H. A CNN Inference Accelerator on FPGA With Compression and Layer-Chaining Techniques for Style Transfer Applications. IEEE Transactions on Circuits and Systems I: Regular Papers, 2023, 70(4): 1591-1604.

[14] Yang Y, Zhu D. Environmental design vs. Environmental art design: a Chinese perspective. Journal of History Culture and Art Research, 2020, 9(4): 122-133.

[15] Wang R. Computer-aided interaction of visual communication technology and art in new media scenes. Computer-Aided Design and Applications, 2021, 19(S3): 75-84.

[16] Liu J, Prouzeau A, Ens B. Effects of Display Layout on Spatial Memory for Immersive Environments. Proceedings of the ACM on Human-Computer Interaction, 2022, 6(ISS): 468-488.

[17] Liu C, Lin M, Rauf H L, Shareef S. Parameter simulation of multidimensional urban landscape design based on nonlinear theory. Nonlinear Engineering, 2022, 10(1): 583-591.

[18] Wang H. Landscape design of coastal area based on virtual reality technology and intelligent algorithm. Journal of Intelligent & Fuzzy Systems, 2019, 37(5): 5955-5963.

[19] Thamrin D, Wardani L K, Sitindjak R H I, Natadjaja L. Experiential learning through community co-design in Interior Design Pedagogy. International Journal of Art & Design Education, 2019, 38(2): 461-477.

[20] Fang Y, Luo B, Zhao T, Jiang B B, Liu Q L. ST-SIGMA: Spatio-temporal semantics and interaction graph aggregation for multi-agent perception and trajectory forecasting. CAAI Transactions on Intelligence Technology, 2022, 7(4): 744-757.

# An Improved K-means Clustering Algorithm Towards an Efficient Educational and Economical Data Modeling

Rabab El Hatimi, Cherifa Fatima Choukhan, Mustapha Esghir

Laboratory of Mathematics, Computing and Applications, Faculty of Sciences,
Mohammed V University in Rabat, Rabat 10000, Morocco

*Abstract*—Education is one of the most crucial pillars for the sustainable development of societies. It is essential for each country to assess its level of access to education. However, the conventional methods of ranking access to education have their limitations. Therefore, there is a need for strategic planning to develop a new classification methods. This study aims to address this need by developing an innovative and efficient unsupervised K-Means model capable of predicting global access to education. The novel approach adopted in this research fills a gap in traditional ranking methods for assessing access to education. Utilizing statistical analysis of data sourced from the World Bank, we evaluated education access across 217 countries spanning various continents and levels of development. By employing economic and educational factors as input for the K-Means algorithm, we successfully identified three distinct clusters, each comprising countries with similar levels of education access. The reliability of our approach was reinforced through rigorous statistical testing to validate the results. Furthermore, we compared the economies of countries within each cluster using primary data, enabling specific recommendations at the economic level to assist countries with limited education access in enhancing their circumstances. Finally, this study makes a significant contribution by introducing a new approach to globally assess education access. The findings provide practical recommendations to aid countries in improving their educational opportunities.

*Keywords—Education assessment; unsupervised learning; statistical analysis; world bank data; K-means*

## I. INTRODUCTION

Education is the foundation of human life and is indispensable for sustainable development. Without education, envisioning a prosperous society is impossible. Science and knowledge have played a crucial role in developing practical applications that meet human needs and improve their quality of life. Education is an essential catalyst in reducing poverty and achieving sustainable development goals. It is a fundamental human right that extends throughout life.

The global state of education, particularly for children, remains a pressing concern in many countries, especially in developing nations. Despite an increase in literacy rates for individuals aged 15 and above since 2000, reaching 85% according to UNESCO statistics, progress made remains fragile and vulnerable to various factors such as dropout rates, livelihood challenges, economic difficulties, global crises, and more [1].

For instance, in 2020, the COVID-19 pandemic resulted in temporary school closures in the majority of countries, affect-ing over 91% of students worldwide. While some developed countries managed to adapt by implementing remote learning or taking necessary measures to resume education, many developing countries faced significant obstacles in resuming normal schooling, leading to increased illiteracy rates and decreased access to education for children [2], [3], [4].

Ranking countries based on the percentage of access to education is closely related to various indicators. It is generally assumed that resource-rich countries such as those with oil, phosphate, and gold would benefit from favorable access to education. However, this assumption is not always valid as other factors come into play, influencing the effective utilization of those resources. Human resources play a crucial role in realizing the full potential of natural resources. This can be illustrated by comparing oil-producing European countries with strong economies and high levels of education to certain African countries that possess significant natural resources but struggle with weak economies and limited access to education [5], [6], [7].

Education is not limited to the availability of educational institutions but encompasses aspects such as the quality of teaching, equal opportunities, inclusion of marginalized groups, and relevance of educational programs. Access to quality education is a fundamental right of every individual, regardless of their socioeconomic background, gender, disability, or place of residence [16].

Analyzing economic data and educational performance of countries can provide valuable insights for policymakers to understand the specific challenges different countries face and develop appropriate education policies. The use of unsupervised learning techniques like clustering can help identify groups of countries that are similar in terms of economic situation and access to education, which can provide benchmarks for effective strategy and policy development [26].

Technically speaking, World Bank data offers a wealth of information on numerous countries worldwide, including indicators such as gross domestic product, domestic product per capita, foreign direct investment, net inflows, balance of payments, inflation, prices, consumption, population, poverty rates, and more [8]. These data can be leveraged through mathematical algorithms to generate statistics that have a positive impact on our world. As a result, numerous studies have utilized World Bank data to make informed decisions on topics such as school dropout rates, student performance, the relationship between education and sustainable development,

and more [9], [10], [11], [12], [13]. The proposed work contributes to the existing body of knowledge in the following ways:

- We developed a novel strategy based on an unsupervised classification model that utilizes World Bank data to rank countries according to their percentage of access to education. The data analysis process involves two steps: statistical analysis and data processing using K-Means with different cluster numbers ($K = 2$ and $K = 3$).

- Significance of statistical tests and selection techniques in the clustering model. Variable selection techniques provide valuable insights into the most important factors, while statistical tests help determine the most efficient clustering results.

- Our study allowed us to determine two risk groups and to identify common problems in these countries to help decision makers to make a good decision.

- This study serves as a stepping stone towards the development of a semi-autonomous and rapid diagnostic system. Such a system would be valuable for predicting access to education in future updates of country data and would align with the objectives of the Sustainable Development Goals, offering substantial opportunities for progress.

The remaining sections of this article are structured as follows: The Materials and Methods in Section III provides an overview of the dataset used, explains the data preprocessing steps, and outlines the clustering algorithm employed. Next, the Results and Discussion in Section IV presents and discusses the findings of the experiment. Finally, the concluding part in Section V offers insights and perspectives on the potential implications of this analysis approach.

## II. RELATED WORK

Several studies have been conducted to examine the links between state regulation in the education sector, access to education, social inequalities, and the use of data in education research.

Vorontsova et al.[14] conducted a comprehensive analysis of the relationship between state regulation indicators in the education sector and the achievement of sustainable development goals. Their study focused on countries in Central and Eastern Europe and used World Bank data from 2006 to 2016. They found that state efforts in the education sector significantly influenced the achievement of sustainable development goals.

Zhongming et al. [15] studied the impact of declining access to education on learning outcomes, particularly in developing countries. Their study, which covered 87 countries and individuals born between 1950 and 2000, revealed that completing at least five years of schooling was crucial for acquiring reading and writing skills. They attributed the disparity in educational quality to social inequalities prevalent in developing countries.

The challenges faced by education in many countries, such as economic difficulties and poverty, have been widely recognized [16], [17]. Social disparities within these nations have

also been identified as factors affecting academic outcomes [18], [19]. Furthermore, the development of education in rural areas is hindered by various developmental factors, including the lack of basic infrastructure [20], [21].

Despite these challenges, some developing countries have made significant progress in the field of education. For example, Rwanda has demonstrated its commitment to education development by allocating a substantial share of state revenues to the sector [22], [23]. Although the country's progress is not without flaws due to political priorities and the transition to a different education system, access to education has significantly improved, and illiteracy rates have decreased [24].

Masino et al. [25] conducted a study on countries that have made notable advancements in education. Their research highlighted factors such as adequate resources, effective policies, community management, decentralization reforms, knowledge dissemination, and increased community participation as key contributors to education development.

In terms of research methodologies, World Bank data has been widely used to inform decision-making in the field of education. Many studies have leveraged this data to explore topics such as dropout rates, student performance, the relationship between education and sustainable development, and more [9], [10], [11], [12], [13] used machine learning models to classify students based on their learning abilities, achieving high accuracy in predicting students' academic outcomes.

Our work builds upon existing research by proposing a new strategy based on unsupervised classification modeling to assess access to education in different countries. Unlike previous approaches, our method utilizes World Bank data and integrates statistical analysis and data preprocessing using the K-Means algorithm with different numbers of clusters (K=2 and K=3). By employing variable selection techniques and statistical tests, we identify the most important factors and achieve more effective clustering results. This approach enables us to identify at-risk groups and highlight common issues countries face in terms of access to education. This information is crucial for policymakers in making informed decisions. Furthermore, our study paves the way for the development of a rapid and semi-autonomous diagnostic system that could predict access to education in future national data updates and contribute to sustainable development goals. In summary, our work brings a fresh perspective and significant opportunities for advancing the field of assessing access to education.

## III. MATERIALS AND METHODS

### A. Dataset

Data collection is a crucial component of any data analysis project. In our study, we recognized the numerous factors that impact the percentage of access to education, including those directly associated with the education sector, such as age of school enrollment (both at the primary and secondary levels), student-to-class ratio, and schools-to-population ratio, etc [26]. Additionally, there are country-level factors to consider, such as unemployment rate, poverty level, gross domestic product, domestic product per capita, foreign direct investment, net inflows, balance of payments, inflation, prices, consumption, and population size, etc [27]. Recognizing the wealth of

relevant data provided by the World Bank, which encompasses the aforementioned indicators, we carefully prepared a comprehensive dataset to extract insightful information and enhance the depth of our study.

*1) Data description:* Initially, we obtained the data from the World Bank website and selected 60 variables that were deemed relevant for our study. These variables encompass information for 217 countries spanning a 20-year period, specifically from 2001 to 2020 (utilizing the latest available statistics for each indicator). As a result, our dataset comprises a comprehensive and representative sample of data, encompassing various country categories, including developed, developing, poor, and marginalized nations.

*2) Data cleaning:* The initial version of the dataset, encompassing all 217 countries, contained numerous missing values (NaN). These missing values indicate either countries withholding their information or having no transactions recorded with the World Bank. To address this, we initiated the data cleaning process by eliminating features with more than 50% missing values. For the remaining missing values, we applied an imputation technique by replacing them with the average value of the respective feature.

Finally, we obtained a thoroughly cleaned and prepared dataset comprising 217 countries and 25 variables. These variables were divided into two categories: 13 variables related to educational factors and 12 variables related to economic factors. The Table I below provides the names of the variables utilized in our study.

TABLE I. Educational and Economical Variables

| Educational variables | Unschooled adolescents; Unschooled kids; Higher education inscriptions; Preschool inscriptions; Primary school inscriptions; Secondary school inscriptions; Ratio female/male in higher education; Ratio girls/boys in primary school; Ratio girls/boys in secondary school; Primary school achievement rate; Secondary school achievement rate; Youth literacy rate; Total literacy rate. |
|---|---|
| Economical variables | Financing capacity; Unemployment; GPD growth; GPD per capita growth; Gini index; Labor force by the level of education; Employment rate 15+; Employment rate 15-24 years; Gross saving; Income share held by highest 20%; Income share held by lowest 20%; Gross domestic saving (% of GPD). |

*3) Data exploration:* In order to make a descriptive analysis and evaluate the dependence of the variables, we constructed two correlation matrix which respectively represent the educational and economic variables (see Fig. 1 and 2).

In Fig. 1, we presented the correlation among the educational variables. The results show a high correlation between all the variables, indicating that the selected variables are reliable and follow a similar distribution. This is a positive indicator for our study. For instance, the variables related to unschooled individuals (Unschooled adolescents, Unschooled



Fig. 1. Correlation matrix for educational variables.



Fig. 2. Correlation matrix for economic variables.

kids, and Higher education enrollments) exhibit strong correlations with all the variables. Additionally, we observe favorable correlations between several variables, such as primary and secondary school achievement rates with school success rates, higher education enrollments, and success in secondary school, among others.

In Fig. 2, we present the correlation between economic variables. Overall, we observe that there is not a strong correlation between the variables, but there are specific variables that demonstrate significant correlations that reveal hidden information. For example, the financing capacity of a country is positively correlated with gross savings. Additionally, variables related to income distribution show a negative correlation with the GINI index, indicating that countries with a more equal

Fig. 3. Eigenvalue scaling



Fig. 4. Circle of correlations of F1 and F2.

income distribution tend to have a lower GINI index. Furthermore, we observe that higher education for development has a positive impact on the employment of the population aged 15 and above. As a result, we attempted to remove some variables that exhibit full correlation since they contain redundant information, such as the Gini index and the share of income held by the highest. It is important to note that a strong correlation between variables is vital for machine learning algorithms to effectively train and demonstrate the quality of the data used.

*4) Feature selection:* To conduct feature selection, we carried out a principal component analysis to identify the variables that will be included in our dataset, as well as the countries that have the most significant contributions to these variable.

*a) Educational data:* According to the Kaiser criterion [29], the inertia for our PCA [30] must be greater than 7.14% (see Fig. 3). We observe that the F1 component explains almost 61% of the variance, and the F2 component explains 10.4%. Our analysis will therefore focus on the F1 and F2 components (see Fig. 4).

- F1: Unschooled adolescents, unschooled kids, preschool inscriptions, ratio female/male in higher education, secondary school inscriptions, primary school achievement rate, secondary school achievement rate, youth literacy rate and total literacy rate.

- F2: Progression to secondary school, higher education inscriptions, primary school inscriptions, ratio female:male in primary school, ratio female:male in secondary school.

*b) Economic data:* According to the Kaiser criterion, components that explain more than 8.33% (see Fig. 5) of



Fig. 5. Eigenvalue scaling.

the variance can be selected. Therefore, we will focus on components F1, F2, F3, and F4 for our analysis. (see Fig. 6).

Fig. 6. Circles of correlations F1, F2 and F3, F4.

- F1: GINI index, income share held by the lowest 20%, income held by the highest 20%, labor force with basic education, population ages 0–14.

- F2: Total unemployment, GDP per capita, employment rate 15⁺, labor force participation rate 15–24, gross domestic savings.

- F3: GDP growth, GDP per capita growth.

- F4: Net lending $(+)$/net borrowing $(-)$ , gross savings (of GDP).

The analysis of economic data does not provide as strong evidence as the education data. The distinction between rich and poor countries is less apparent, and it seems that economic factors are not the sole determinants of non-schooling in children. Therefore, classification methods can be used to group countries with similar characteristics. This allows for more targeted strategies to improve the educational performance of children and adolescents.

### B. Classification of Countries according to Access to Education

Unsupervised learning is a branch of machine learning that involves analyzing and grouping unlabeled data [28]. These algorithms aim to identify patterns or clusters in the data without the need for explicit human guidance. In mathematical terms, unsupervised learning involves observing multiple occurrences of a vector X and learning the probability distribution p(X) for those occurrences. One of the most commonly used algorithms in recent years for unsupervised learning is K-means [31] [32], [33], [34].This method partitions data points into k clusters, where each data point is assigned to the cluster that is closest to it. The objective function for this method is to sum the squared distances within each cluster, across all clusters:

$$\underset{S}{\operatorname{argmin}} \sum_{i=1}^{k} \left( \sum_{x_j \in S_i} ||x_j - \mu_i||^2 \right)$$

where:
$x_j$ is a data point in the dataset
$S_i$ is a cluster (collection of data points)
$\mu_i$ is the cluster mean (the centre of cluster $S_i$).

In our case, after preparing the data and selecting the important features, we tried to apply the k-mean algorithm to cluster these data. The purpose of this analysis is to assign labels to our dataset based on the level of access to education in order to achieve the most optimal classification of countries.

Fig. 7 illustrates the approach followed in this study.

### IV. RESULTS AND DISCUSSION

#### A. Default Number of Clusters

The k-means method requires us to specify the number of clusters we want. To determine the optimal number of clusters for our dataset, we used the elbow technique [35],which considers the distortions of the variables in our study. By analyzing the graph, we observed that the points after which the distortion starts to decrease are 2 and 3. Therefore, the optimal number of clusters that can be used for our data is either 2 or 3 clusters (see Fig. 8).

#### B. K-means Clusters

*a)* ***For*** $k = 2$*:* We initially tested the case with $k = 2$ to determine if we could identify nations with high unschooling rates (see Fig. 9) before proceeding with the k-means analysis using $k = 3$. The results of the two clusters are presented on a world map to aid in interpretation.

The results of the two clusters are plotted on a world map to facilitate interpretation (see Fig. 10).

Fig. 7. Implementation of $k$-means clustering on our dataset.



Fig. 8. Elbow method on distortions.



Fig. 9. Abstract representation of the k-means clusters with $k = 2$.

As observed on the map (see Fig. 10), the initial results indicate that Central African countries and some Central Asian countries form a distinct cluster representing nations with limited access to education. On the other hand, the majority of countries from other continents are grouped together in the cluster representing countries with better access to education. While these findings are expected and widely known, our objective is to identify countries that are making efforts to improve their education systems and analyze those with the highest levels of access to education. Therefore, we will explore the option of increasing the number of clusters to achieve a more refined classification.

*b) For* $k = 3$: To further our interpretation and understanding of differences in access to education, we now consider the clustering when $k = 3$ (see Fig. 11).
A geographic mapping is shown in Fig. 12:

*C. Validation of Results Using Statistical Analysis*

In order to validate the obtained results for both $k = 2$ and $k = 3$, we conducted a statistical analysis. The analysis involved comparing the cluster means to determine the most effective clustering result (see Table II and Table III). To facilitate this comparison, we utilized a diagram, as shown in Fig. 13.

To assess the equality of population variances, we employed Bartlett's test. The null hypothesis $H_0$ states that all k population variances are equal, while the alternative hypothesis $H_1$ suggests that at least two variances are different. If the p-value is greater than $0.05$, the null hypothesis is retained,

Fig. 10. Map representing the countries belonging to the clusters with $k = 2$.



Fig. 12. Map representing the countries belonging to the clusters with $k = 3$.



Fig. 11. Clusters with $k = 3$.



Fig. 13. Methodology for comparing averages.

indicating that the samples are identical. Conversely, if the p-value is less than $0.05$, we reject the null hypothesis and apply Welch's test to estimate the averages [36].

- **k-Means with $k = 2$**

The statistical analyses reveal significant differences in the mean scores of education-related variables between the two clusters. Cluster $0$ is characterized by low non-enrollment rates, high enrollment rates across all education levels, higher completion rates, and high literacy rates, as shown in Fig. 14. In terms of gender ratios, we observe that countries in Cluster 1 have lower female-to-male ratios compared to Cluster 0 (see Fig. 15). However, the two clusters do not exhibit differentiation based on any variables related to economic health, except for factors related to the active population, unemployment, and gross savings (see Fig. 16).

- **Clustering with $k = 3$**

The statistical tests indicate that the second test provides better discrimination for our sample. The three clusters show significant differences in terms of educational variables. Clusters 0 and 2 exhibit higher levels of non-schooling among children and adolescents, lower enrollment rates in secondary and higher education, and lower literacy rates, as depicted in Fig. 17. Furthermore, the gender ratios in Cluster 2 are statistically lower compared to Clusters 0 and 1 (see Fig. 18).

On the economic front, Clusters 1 and 2 stand out for their lower unemployment rates, significant GDP growth (both total and per capita), and higher rates of young individuals in the labor force, indicating greater financing needs (see Fig. 19).

We observe that the average scores of the variables related to education are similar across the three clusters. However, the economic data are able to discriminate among the clusters (see Fig. 17 18 19). This suggests that economic factors played a significant role in the final grouping ($k = 3$), contributing substantially to the classification of countries based on access to education. The resulting clustering appears to be realistic and acceptable.

Fig. 14. Average scores of different education-related variables in our two clusters.



Fig. 15. Average parity ratios in the different school cycles in our two clusters.



Fig. 16. Average scores of economic variables in our two clusters.



Fig. 17. Average scores of educational variables in the three clusters.



Fig. 18. Average parity ratios in the different school cycles in the three clusters.



Fig. 19. Averages of economic variables in the three clusters.

TABLE II. STATISTICAL TESTS FOR THE EDUCATIONAL VARIABLES

| Variables | Bartlett's test | Student's test | Welsh's test |
|---|---|---|---|
| unschooled_teenager | $p = 7.11e - 21$ | | $p = 7.12 - 12$ |
| unschooled_kids | $p = 2.46e - 20$ | | $p = 6.45e - 09$ |
| higher_educ_regis | $p = 1.00e - 09$ | | $p = 4.95e - 17$ |
| primary_regis | $p = 0.00$ | | $p = 0.01$ |
| preschool_regis | $p = 0.04$ | | $p = 1.12e - 13$ |
| secondary_regis | $p = 0.32$ | | $p = 4.32e - 32$ |
| ratio_higher_educ | $p = 0.27$ | $p = 2.28e - 17$ | |
| ratio_primary | $p = 3.83e - 26$ | | $p = 7.32e - 06$ |
| ratio_secondary | $p = 5.53e - 16$ | | $p = 1.31e - 09$ |
| rate_prim_comp | $p = 1.23e - 08$ | | $p = 1.26e - 18$ |
| rate_second_comp | $p = 0.49$ | | $p = 6.14e - 31$ |
| tee_literacy_rate | $p = 5.86e - 21$ | | $p = 1.61e - 12$ |
| tee_literacy_rate | $p = 5.02e - 07$ | | $p = 2.90e - 15$ |

TABLE III. STATISTICAL TESTS FOR THE ECONOMIC VARIABLES

| Variables | Bartlett's test | Student's test | Welsh's test |
|---|---|---|---|
| financing_capacity | $p = 2.29e - 24$ | | $p = 0.26$ |
| total_unemployement | $p = 0.62$ | $p = 0.006$ | |
| gip_increasing | $p = 0.08$ | $p = 0.04$ | |
| gip_increasing_inhab | $p = 0.05$ | $p = 0.70$ | |
| index_Gini | $p = 0.06$ | $p = 0.98$ | |
| active_graduate_pop | $p = 0.10$ | $p = 0.01$ | |
| employ_rate_15$^+$ | $p = 0.22$ | $p = 0.004$ | |
| active_rate_15–24 | $p = 0.003$ | | $p = 0.007$ |
| gross_saving (GIP) | $p = 1.78e - 07$ | $p = 0.02$ | |

In conclusion, it can be inferred that the economy of a country is closely linked to education. Countries with stronger economies tend to prioritize education.

### D. Discussion of Results

The complete World Bank dataset comprises multiple years of information (2001 to 2020) for most countries. This dataset was utilized for statistical analysis, specifically for developing k-means clustering models. The objective was to predict the level of education access and the economic health of the countries in our sample. The analysis focused solely on standardized data related to education access and the economic conditions of the countries.

The study began with an assessment of the overall state of education, which revealed positive trends in non-enrollment, access, and completion rates, except for several African countries with lower levels of education access. Economic data highlighted concerns regarding wealth distribution in certain countries. Subsequently, an unsupervised machine learning model, specifically the k-means algorithm, was employed to segment the data into two cases (k=2 and 3). The k=2 test aimed to identify clear distinctions between northern and southern countries. The results demonstrated a significant disparity, indicating that developed countries have lower levels of illiteracy compared to other countries, which appear isolated in terms of education access.

When setting $k = 3$, we observe that the educational data provide less discrimination compared to the $k = 2$. However, based on the statistical analysis of the tests, the model with $k = 3$ performs better overall. With this approach, we obtain three distinct clusters that exhibit significant differences in various educational factors (see Fig. 17 and Fig. 18), and economic variables (see Fig. 19 and 12), including economic

growth and labor force. By segmenting the data, our algorithm enables the identification of two risk groups and facilitates the recognition of common challenges faced by these countries. Cluster 1 is composed essentially of countries in Africa and in Arabian Peninsula like Chad, Yemen, Syria, and Afghanistan. These countries have high rates of non-enrollment but tend to be closer to those of clusters 0 and 2 17. However, their financing needs are still significant. And as a result, these countries where the action is necessary are also different from those of cluster 2 19. We note that these countries experienced civil wars starting in the 2000. In these countries, schooling is compulsory until the age of 14 on average. After this age, children work under challenging circumstances, which explains the low unemployment rates and the high rates of non-enrollment among adolescents.

Cluster 2 countries have the highest out-of-school rates with high financing needs 17 19. Nevertheless, they also have the lowest unemployment rates for all age groups. Wealth is still evenly distributed. Therefore, these countries require special attention. Among them are Zimbabwe, Papua New Guinea, Somalia, Morocco, Kenya, India, and Suriname. These nations are former colonies and must now adapt to independence. They display high rates of inequality in access to education due to different social systems (castes in India, urban vs rural children in Somalia or Zimbabwe, etc.) In most of these countries, schooling is not a free choice, which does not give equal access to all. These countries are also marked by numerous conflicts (Papua New Guinea, Somalia), where children are enrolled very early in the armed conflict. Finally, there are significant disparities in access based on gender. Girls are more likely to be sold, mutilated, or forced into marriage. Therefore, it would seem that accompanying measures are the most effective. Indeed, the countries in question are developing countries that find themselves having to manage their independence acquired some 40 years ago. The years of colonization did not prepare them to handle their economy in a way that would have allowed them to know where to invest their funds. These are resource-rich countries that need guidance and human assistance to progress.

## V. CONCLUSION

Our statistical analysis and utilization of machine learning have highlighted the complexity of the global education problem, calling for a comprehensive approach from organizations tasked with improving access to education. It is crucial for countries facing conflicts, whether civil or otherwise, to put an end to these situations in order to restore and uphold children's rights to education.

Furthermore, we have identified inequalities in access to education, particularly based on gender and geographic origin of children. For instance, children living in war-torn areas face significant barriers to education. As we continue our research, it would be interesting to expand our model to encompass other factors beyond just the economy and education. While we utilized economic and educational variables, our findings were enriched by in-depth field studies, enabling us to better understand countries grappling with education issues and the underlying economic causes.

This project carries significant implications for organizations such as the Council of Europe, UNICEF, and UNESCO,

which are committed to promoting education for all worldwide. Governments of individual countries can also find valuable guidance for enhancing access to education within their own borders.

In conclusion, our research underscores the need for concrete policies and actions to improve access to education in countries with the lowest levels of access. By combining data-driven approaches with on-the-ground studies, we can make meaningful contributions toward the shared goal of quality education for all children worldwide.

For future work, it would be valuable to further explore the impact of social and cultural factors on access to education. By incorporating these additional variables into our model, we could gain a deeper understanding of the specific challenges faced by certain demographic groups or regions (especially war-torn countries).

Additionally, it would be beneficial to examine technology-based solutions for enhancing access to education. Leveraging technology can help overcome geographical barriers and provide educational resources to remote or underdeveloped regions. Studies on the effectiveness and acceptance of these technological solutions in different contexts could be conducted.

Lastly, it would be relevant to investigate the long-term impact of improving access to education on a country's economic and social development. By evaluating medium and long-term outcomes, we could gain a better understanding of the benefits and ripple effects of investment in education.

These future endeavors could contribute to strengthening our understanding of access to education and inform policies and actions aimed at promoting inclusive and quality education for all.

## REFERENCES

[1] Fernández Álvarez, R. (2020). Geoparks and education: UNESCO Global Geopark Villuercas-Ibores-Jara as a case study in Spain. Geosciences, 10(1), 27.

[2] Rao, N., & Fisher, P. A. (2021). The impact of the COVID-19 pandemic on child and adolescent development around the world. Child development, 92(5), e738.

[3] Memon, A. S., Rigole, A., Nakashian, T. V., Taulo, W. G., Chávez, C., & Mizunoya, S. (2020). COVID-19: How prepared are global education systems for future crises?.

[4] Munir, F. (2021). Mitigating COVID: Impact of COVID-19 lockdown and school closure on children's well-being. Social Sciences, 10(10), 387.

[5] Cole, J. (2021). PL5. Training and education in Europe, Middle East, Africa, Latin America and Asia Oceania chapters, IFCN: An international survey. Clinical Neurophysiology, 132(8), e38.

[6] de Jong, W., Huang, K., Zhuo, Y., Kleine, M., Wang, G., Liu, W., & Xu, G. (2021). A Comparison of Forestry Continuing Education Academic Degree Programs. Forests, 12(7), 824.

[7] Nguyen, T. P. L., Nguyen, T. H., & Tran, T. K. (2020). STEM education in secondary schools: Teachers' perspective towards sustainable development. Sustainability, 12(21), 8865.

[8] Dang, H. A. H., Pullinger, J., Serajuddin, U., & Stacy, B. (2021). Statistical Performance Indicators and Index.

[9] Chai, K. E., & Gibson, D. (2015). Predicting the Risk of Attrition for Undergraduate Students with Time Based Modelling. International Association for Development of the Information Society.

[10] Mundy, K., & Verger, A. (2016). The World Bank and the global governance of education in a changing world order. The handbook of global education policy, 335-356.

[11] Vladimirova, K., & Le Blanc, D. (2016). Exploring links between education and sustainable development goals through the lens of UN flagship reports. Sustainable Development, 24(4), 254-271.

[12] Vladimirova, K., & Le Blanc, D. (2016). Exploring links between education and sustainable development goals through the lens of UN flagship reports. Sustainable Development, 24(4), 254-271.

[13] Zhang, T., Shaikh, Z. A., Yumashev, A. V., & Chłąd, M. (2020). Applied model of E-learning in the framework of education for sustainable development. Sustainability, 12(16), 6420.

[14] Vorontsova, A. S., Vasylieva, T. A., Bilan, Y. V., Ostasz, G., & Mayboroda, T. (2020). The influence of state regulation of education for achieving the sustainable development goals: case study of Central and Eastern European countries.

[15] Le Nestour, A., Moscoviz, L., & Sandefur, J. (2022). The long-run decline of education quality in the developing world. Center for Global Development.

[16] Mansi, E., Hysa, E., Panait, M., & Voica, M. C. (2020). Poverty—A challenge for economic development? Evidences from Western Balkan countries and the European Union. Sustainability, 12(18), 7754.

[17] Maneejuk, P., & Yamaka, W. (2021). The impact of higher education on economic growth in ASEAN-5 countries. Sustainability, 13(2), 520.

[18] Emmerich, M., & Hormel, U. (2021). Unequal Inclusion: The production of social differences in education systems. Social Inclusion, 9(3), 301-312.

[19] Thoman, D. B., Lee, G. A., Zambrano, J., Geerling, D. M., Smith, J. L., & Sansone, C. (2019). Social influences of interest: Conceptualizing group differences in education through a self-regulation of motivation model. Group Processes & Intergroup Relations, 22(3), 330-355.

[20] Fafunwa, A. B., & Aisiku, J. U. (Eds.). (2022). Education in Africa: A comparative survey. Taylor & Francis.

[21] Solstad, K. J., & Karlberg-Granlund, G. (2020). Rural education in a globalized world. Educational Research and Schooling in Rural Europe: An Engagement with Changing Patterns of Education, Space and Place. Charlotte: Information Age Publishing, 49-76.

[22] Williams, T. (2016). Oriented towards action: The political economy of primary education in Rwanda. Effective States and Inclusive Development (ESID) Working Paper, (64).

[23] Mann, L.,& Berry, M. (2016). Understanding the political motivations that shape Rwanda's emergent developmental state. New Political Economy, 21(1), 119-144.

[24] Williams, T. P. (2017). The political economy of primary education: Lessons from Rwanda. World Development, 96, 550-561.

[25] Masino, S., & Niño-Zarazúa, M. (2016). What works to improve students' learning outcomes in developing countries? Revue internationale de développement de l'éducation, 48, 53-65.

[26] Morgan, T. L., Zakhem, D., & Cooper, W. L. (2018). From high school access to postsecondary success: An exploratory study of the impact of high-rigor coursework. Education Sciences, 8(4), 191.

[27] Mia, M. M., Zayed, N. M., Islam, K. M. A., Nitsenko, V., Matusevych, T., & Mordous, I. (2022). The Strategy of Factors Influencing Learning Satisfaction Explored by First and Second-Order Structural Equation Modeling (SEM). Inventions, 7(3), 59.

[28] Shamsuddin, M. R., Abdul-Rahman, S., & Mohamed, A. (2019). Exploratory analysis of MNIST handwritten digit for machine learning modelling. In Soft Computing in Data Science: 4th International Conference, SCDS 2018, Bangkok, Thailand, August 15-16, 2018, Proceedings 4 (pp. 134-145). Springer Singapore.

[29] Cox, D. R., & Lewis, P. A. (1966). The statistical analysis of series of events.

[30] Severson, K. A., Molaro, M. C., & Braatz, R. D. (2017). Principal component analysis of process datasets with missing values. Processes, 5(3), 38.

[31] Ahmed, M., Seraj, R., & Islam, S. M. S. (2020). The k-means algorithm: A comprehensive survey and performance evaluation. Electronics, 9(8), 1295.

[32] Moussaid, A., Fkihi, S. E., & Zennayi, Y. (2021). Tree crowns segmentation and classification in overlapping orchards based on satellite images and unsupervised learning algorithms. Journal of Imaging, 7(11), 241.

[33] Sinaga, K. P., & Yang, M. S. (2020). Unsupervised K-means clustering algorithm. IEEE access, 8, 80716-80727.

[34] Yuan, C., & Yang, H. (2019). Research on K-value selection method of K-means clustering algorithm. J, 2(2), 226-235.

[35] Bholowalia, P., & Kumar, A. (2014). EBK-means: A clustering technique based on elbow method and k-means in WSN. International Journal of Computer Applications, 105(9).

[36] Marco-Franco, J. E., Reis-Santos, M., Barrachina-Martínez, I., González-de-Julián, S., & Camaño-Puig, R. (2022). Validation of a new telenursing questionnaire: testing the test. Mathematics, 10(14), 2463.

# Investigating the Impact of Preprocessing Techniques and Representation Models on Arabic Text Classification using Machine Learning

Mahmoud Masadeh[1][¶], Moustapha.A[2][¶], Sharada B[3], Hanumanthappa J.[4], Hemachandran K[5],
Channabasava Chola[*][6], Abdullah Y. Muaad[*][7]

Computer Engineering Department, Yarmouk University, Irbid 21163, Jordan[1]
Department of Studies in Computer Science, Mysore University, Manasagangothri, Mysore 570006, India[2,3,4,7]
Department of Business Analytics, School of Business, Woxsen University, Hyderabad, India 502345[5]
Department of Electronics and Information Convergence Engineering, College of Electronics and Information,
Kyung Hee University, Suwon-si 17104, Republic of Korea[6]
[¶]The authors have an equal contribution

*Abstract*—Arabic Text Classification (ATC) is a crucial step for various Natural Language Processing (NLP) applications. It emerged as a response to the exponential growth of online content like social posts and review comments. In this study, *preprocessing techniques* and *representation models* are used to evaluate the effectiveness of ATC using Machine Learning (ML). Generally, the ATC operation depends on various factors, such as stemming in preprocessing, feature extraction and selection, and the nature of the dataset. To enhance the overall classification performance, preprocessing methodologies are primarily employed to transform each Arabic term into its root form and reduce the dimensionality of representation. In the representation of Arabic text, feature extraction and selection processes are imperative, as they significantly enhance the performance of ATC. This study implements the chosen classifiers using various feature selection algorithms. The comprehensive assessment of classification outcomes is conducted by comparing various classifiers, including Multinomial Naive Bayes (MNB), Bernoulli Naive Bayes (BNB), Stochastic Gradient Descent (SGD), Support Vector Classifier (SVC), Logistic Regression (LR), and linear Support Vector Classifier (LSVC). These ML classifiers are assessed utilizing short and long Arabic text benchmark datasets called BBC Arabic corpus and the COVID-19 dataset. The assessment findings indicate that the efficacy of classification is significantly influenced by the preprocessing methods, representation model, classification algorithm, and the datasets' characteristics. In most cases, the SGDC and LSVC have consistently surpassed other classifiers for the datasets under consideration when significant features are chosen.

*Keywords*—*Arabic Text Classification (ATC); Text Mining (TM); Machine Learning (ML); preprocessing methods; representation models; Feature Extraction (FE); Feature Selection (FS)*

## I. INTRODUCTION

Text Mining (TM), i.e., knowledge discovery, which entails extracting meaningful information from text, has gained significant attention in recent years. With the exponentially generated textual data on many social media sites, understanding and analyzing text data become increasingly complex [1]. TM is a discipline that requests the support of other scientific subjects such as statistics, machine learning (ML), natural language processing (NLP), and linear algebra [2].

One of the vital topics in TM is text classification (TC), i.e., categorization, which is a challenging computational task in the TM field [3]. TC has become a big data problem as textual data's volume, variety, and velocity increase rapidly. Furthermore, only a few models and tools can help understand and classify Arabic Text (AT). Therefore, it is compulsory to design efficient models to address AT's challenges and support decision-makers in making the right decisions in many real-life domains such as healthcare, economics, social media, and financial markets [4].

Over 477 million people speak Arabic as their first language. Furthermore, a significant percentage speak Arabic as a second language [5]. Arabic is the official language in 22 nations and the original script for Persian and Urdu [6]. The Arabic language consists of 29 letters written from right to left. The Arabic language distinguishes itself by position-dependent letter forms and shapes [7]. Arabic, unlike other languages, is founded on roots, which contributes to its complexity. Aside from that, Arabic is divided into three formats: Classical Arabic (CA), Modern Standard Arabic (MSA), and Dialectal Arabic (DA). Other AT-related issues include phonology, orthography, and morphology. So, working with AT is more complicated than other languages, such as English [7]. Thus, there is an urgent need to handle Arabic Text Classification (ATC).

Arabic Text Classification starts with preprocessing to make the text ready for further processing. Then, a group of *representation learning methods* are applied. They make the text understandable by machines and automatically find the patterns that will help to discriminate classes and achieve a classification task [2]. Once the representation for a given text collection is created, an optimal set of features is required. Therefore, Feature Extraction (FE) and Feature Selection (FS) techniques must be applied to extract and select the best features to decrease the dimension of the representation process. Subsequently, the classifier is trained to learn the pattern in the training phase (offline learning) and classify the text into different classes in the testing phase (online learning) [8].

---

*Corresponding authors

Text Classification (TC) is a process in which the system will assign a suitable label for each test document based on recognizing what has already been learned from the training phase [9]. The complexity of textual data made classification a challenging task. Therefore, over the past few decades, TC have been extensively researched and addressed in different applications such as mail spam filtering, document classification, web searching, and web page classification. However, today's TC information has become challenging and more significant, mainly with detection tasks such as hate speech and rumours detection, especially in the Arabic language [10].

This study examines the influence of preprocessing techniques and representation models on Arabic text classification, where six ML-based classifiers are utilized for the classification task. Therefore, this research allows developers in the profession to choose a powerful ML-based approach for robust ATC applications. The primary contributions of this study are outlined as follows:

- We investigate the effect of preprocessing, representation and feature selection techniques on Arabic Text Classification (ATC).

- We employ and evaluate many classification models with various preprocessing algorithms to ascertain their efficacy in ATC.

- We demonstrate the efficacy of preprocessing and representation methods; all ML-based classifiers are assessed using two benchmark datasets for Arabic text, specifically BBC Arabic and COVID-19.

- We verify the effectiveness of both long and short Arabic text datasets for ATC tasks.

The rest of this article is presented in the following manner. Section II explains preliminaries essential to comprehend the context of the work. Section III exposes some literature review on preprocessing and representation for ATC. Section IV demonstrates the proposed model of ATC strategy. SectionV explores experimental results, examines a comparative analysis, and deliberates on the outcomes of the experiments. The conclusion of this work is presented in Section VII.

## II. Preliminaries

### A. Natural Language Processing (NLP)

It is an important aspect of Artificial Intelligence (AI). It allows programs and machines to analyze and understand human language, allowing them to carry out repetitive exercises without human intervention. Several domains have developed the basics of NLP, such as artificial intelligence, computer and information sciences, linguistics, electronic and electrical engineering, robotics, math, and psychology. Appliances can analyze and learn human language through a technique known as NLP [11]. NLP-based techniques manipulate a substantial portion of data to fetch valuable information/knowledge. Therefore, various data mining and machine learning strategies are used. Thus, text pre-processing should be involved to prepare the text for other processing, e.g., representation features engineering that is mandated to extract features and hand it to ML techniques [12]. For example, pre-processing could incorporate text tokenization and stop-word removal. Recently,

Arabic NLP has emerged as a nascent research domain. It encompasses the evolution of approaches and tools using the Arabic language. However, it faces multiple complex concerns associated with the form and nature of the Arabic language.

### B. Machine Learning (ML) Algorithms

ML is integrated into various areas, including healthcare [13], hardware design [14] [15], quality control [16], and NLP [12], with this study focusing on the latter. Information is a systematic collection of discrete facts that suppresses the complete range of typical patterns. The primary purpose of machines is to find rituals that remind them of a specific occasion. If the system recognizes these patterns, machine learning has occurred. The authors of [17] stated that advances in machine learning, particularly deep learning, allow us to develop algorithms that utilize real-world data to produce conclusions that look subjective. There are several approaches for preparing text for subsequent processing, as demonstrated in Section IV-C. Text tokenization, also known as text segmentation or linguistic analysis, divides the text into smaller units called tokens, which can be words, characters, or subwords. The most typical method of creating tokens is based on space. Articles (e.g., a, an, the), conjunctions (e.g., and, but, if), and prepositions (e.g., in, at, on) [18] are stop words that do not communicate a clear meaning. As a result, they must be removed. In ML, features are numerical properties. However, in certain cases, such as sentiment analysis, the data may not include numerical qualities. Thus, many forms of features (e.g., word and character) are translated into numerical features, and selecting from them to make ML operate adequately is referred to as feature engineering (*feature selection and feature extraction*).

### C. Text Classification

The number of available complex text documents and the size of the text have just grown exponentially. This mandates a more profound understanding of machine learning techniques to accurately categorize texts in various applications. ML models achieve successful results in NLP since they rely on their ability to comprehend complicated prototypes and non-linear associations within data.

The authors of [19] evaluated the effect of the preprocessing schemes on classification success, considering e-mail and news domains, for Turkish and English. They discovered that selecting suitable combinations of preprocessing tasks, rather than including or excluding them all, may supply substantial enhancement in classification accuracy depending on the target domain and language. The authors of [20] discussed a short recap of text classification algorithms including text feature extractions, dimensionality compaction techniques, existing algorithms, and evaluation methods. The authors of [21] indicated that DL–based models have exceeded classical ML-based techniques in different text classification studies. They introduced an exhaustive study of more than 150 DL–based models for text classification and examined their technical contributions, similarities, and strengths.

### D. Evaluation Metrics

The efficiency of the proposed models is assessed in terms of accuracy, precision, recall, and F1-Score. As presented in

Eq. 1, **accuracy** represents the number of correctly categorized data instances over the whole number of data samples. For an unbalanced dataset, the positive and negative classes have a diverse number of samples. Thus, the accuracy is not appropriate for assessing the model and other metrics are required. **Precision**, i.e., positive predictive value, describes how many of the precisely foreseen cases were positive. As represented in Eq. 2, it should be 1 for an ideal classifier where FP is zero. **Recall**, i.e., sensitivity or true positive rate, is represented in Eq. 3.

Recall for a label is represented as the number of true positives divided by the total number of confirmed positives. For an excellent classifier, recall should be 1 where FN is zero. For a perfect classifier, both precision and recall are 1. **F1-score** is a measure that relies on both precision and recall as represented in Eq. 4. F1-score is 1 when both precision and recall are 1. So, the F1-score is the harmonic mean of precision and recall and it is a more reasonable measure than accuracy [22].

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \qquad (1)$$

$$Precision = \frac{TP}{TP + FP} \qquad (2)$$

$$Recall/Sensitivity = \frac{TP}{TP + FN} \qquad (3)$$

$$F1\ Score = 2 \times \frac{Precision \times Recall}{Precision + Recall} \qquad (4)$$

The above equations can be used with only the binary classification problem. For Multi-Class, the earlier formulas for Precision and Recall might not seamlessly apply. We have to calculate the per-class values of precision, recall, and F1 score, where we have seven classes for the BBC dataset and three classes for the COVID-19 dataset. However, rather than having multiple per-class precision, recall, and F1-score, it would be more suitable to average them to get a single number to represent overall performance. For that, we will use macro average and weighted average.

*Macro averaging* is the most direct among the considerable averaging methods. The macro-averaged of a specific score is calculated by taking the arithmetic mean of all the per-class scores. Thus, treats all classes equally regardless of their support values. Support indicates the number of real occurrences of the class in the dataset. For example, the BBC dataset contains 2,356 documents of class *Middle East News* while containing only 49 documents of class *International Press*. The *weighted-average* of a specific metric, e.g., F1-score, is computed by considering the mean of all per-class F1 scores while taking each class's support. The *weight* indicates the proportion of each class's contribution comparable to the sum of all support values.

## III. LITERATURE REVIEW

This section emphasizes related works on preprocessing and representation in Arabic text classification using machine learning classifiers. The first step in developing ATCs is preprocessing, which includes cleaning text, simplifying computations, and preparing the dataset for further processing [10] [23]. Text preprocessing is essential for preparing it for representation, given that text data is inherently unorganized. Such unstructured text is not amenable to subsequent analysis without a pre-established data model. As a result, specialized preprocessing techniques are necessary to condition the text for further use. Reducing the size of the text and removing extraneous elements could either enhance or detract from its performance [24]. As a second step, converting unstructured text into structured data is called the representation process since machines can only process structured information, whereas AT is unstructured [9] [2].

Reducing the dimensionality of the data is not a compulsory step in ATC but it's often essential. This is due to the enormous number of text features that can generate large and sparse matrices that adversely affect the efficiency of ATCs. Some of these features may add noise, making it difficult to classify different categories accurately. Various methods can be used at this stage, such as Feature Extraction (FE), Feature Selection (FS), and Principal Component Analysis (PCA), among other multivariate techniques that improve performance. Dimensionality can also be lowered in preprocessing through approaches like stemming and lemmatization. It's worth noting that many scholars have used Term Frequency-Inverse Document Frequency (TF-IDF) and Bag-of-Words (BoW) methods, which have three main drawbacks: they create sparse matrices, lose the sequence of words, and neglect semantic context. Thus, leading to identical representations for different sentences. To overcome these constraints, various feature selection methods like Information Gain (IG), Mutual Information (MI), and Chi-square (CHI) have been suggested to identify the best features and reduce dimensionality [10] [25].

## IV. PROPOSED METHODOLOGY

Fig. 1 introduces an Arabic text classification framework using varied preprocessing and representation methods like BoW and TF-IDF. Diverse classifiers are applied, including Multinomial NB, Bernoulli NB, LR, Stochastic Gradient Descent (SGD) Classifier, SVC, and Linear SVC. The study's outcomes highlight the impact of these techniques on enhancing ATC system performance. As depicted in Fig. 1, the system's workflow involves pre-processing, representation, feature selection, and classification algorithms. Initial steps include stop word removal, normalization, and stemming to reduce dimensions. The resulting text undergoes TF-IDF and BoW processes, producing a matrix as input. Machine learning employs this matrix, where the data is divided into 80% for training and 20% for testing.

### A. The Proposed Architecture

The proposed method is divided into four phases. In the first stage, dataset preparation and splitting are carried out. In these data, preprocessing steps, such as normalization,

Fig. 1. Schematic representation of the proposed method.

scaling, and missing values handling, are applied to ensure the data can be used for machine learning. Representation is the second phase, where raw data is transformed into a format that ML algorithms can process. Then, classification, several classifiers are used to categorize text into one or more different classes. Finally, assessment metrics are used to examine the performance and efficacy of a statistical or machine-learning model using quantitative measurements. Fig. 1 depicts the suggested method's design.

### B. DataSet

We selected two Arabic datasets, BBC Arabic and COVID-19, to experimentally evaluate six classifiers based on preprocessing and representation. Performance metrics are analyzed to assess classifier effectiveness [26]. The Arabic BBC corpus dataset comprises 4,763 Arabic documents classified into seven classes. Document distribution per class is Middle East News (2,356), world news (1,489), business and economics (296), sports (219), international press (49), science and technology (232), and art and culture (122) [26]. The corpus contains around 1.8 million words and 106,733 distinct keywords, summarised in Table I. Arabic comments related to COVID-19 are classified using an additional dataset of short texts. These comments are analyzed, and the data distribution of this COVID-19 dataset is presented in Table II which summarizes the details of each benchmark dataset. These datasets belong to two applications (long text and short text). We divided each dataset into training and validation test sets with a ratio of 80% and 20%, according to the Pareto principle [27].

TABLE I. BBC DATASET DISTRIBUTION

| No. | Class Type | No. of Documents |
|---|---|---|
| 1 | Middle East News | 2,356 |
| 2 | Science and Technology | 232 |
| 3 | International Press | 49 |
| 4 | Art and Culture | 122 |
| 5 | Sports | 219 |
| 6 | Business and Economy | 296 |
| 7 | World News | 1,489 |

TABLE II. COVID-19 DATASET DISTRIBUTION

| No. | Class Type | No. of Documents |
|---|---|---|
| 1 | Positive | 7,962 |
| 2 | Negative | 635 |
| 3 | Natural | 1,391 |

### C. Preprocessing

Preprocessing involves transforming raw data into a suitable input format for ML models. It qualifies text by converting it into a convenient form for document classification, reducing complexity. This phase removes non-significant characters, stop words, and punctuation in Arabic text. Preprocessing steps include tokenization, normalization, stop word removal, and stemming.

*1) Tokenization:* Tokenization involves breaking text into tokens and converting words into numbers. The acquired segments can be single items (1-gram) or a sequence of n words (n-gram). In Arabic, sentences are divided using signals like commas, quotes, and spaces. Tokens can be individual words, irrespective of meaning or relationships

*2) Normalization:* Normalization involves converting text letters to a canonical form or removing diacritics, such as changing ة ت ـهـ [23].

*3) Stop Word Elimination:* For a work that targets Arabic text, the first pre-processing step is the elimination of non-Arabic text. Thus, a whitespace character is used instead of each non-Arabic character. Then, we deduct *stopwords*, which occur often in the text and are insignificant for text classification, e.g., هناك ، على ، منذ ، أو ، هو. A list of the most often used Arabic stop words is accessible at [18]. Stopwords account for around 20%-30% of a document's exhaustive words. These terms can be deleted since they are repetitive [28]. The basic approach for extracting stopwords is static, meaning it uses a pre-filled list of all words that are semantically irrelevant to a specific language. For the dynamic approach, stopwords are recognized online rather than previously established, and attributes are given depending on their relevance. In this effort, analogous to the removal of

stopwords, we removed punctuation and numerals from the Arabic text.

*4) Stemming:* Stemming involves removing prefixes, suffixes, and definite articles from words to reach their root form. This process, like root, light, and hybrid stemming, aims to simplify words for analysis. As in the shown example, the words in the sentence are given in their root form.

### D. Text Representation/Feature Engineering

Text representation involves transforming unstructured text into manageable forms. Various text feature representation methods exist, each with distinct traits. The initial dataset includes a group of documents with many classes as expressed in Eq. 5. The frequency (TF) representation technique is the average occurrence within a specific topic divided by the occurrence rate.

$$AD = d_1, d_2, d_3, ......, d_n \qquad (5)$$

$$tf(t, d) = \frac{f_d(t)}{Max_{w \in d} \ f_d(w)} \qquad (6)$$

$$idf(t, D) = Ln[\frac{|D|}{|\{d \in D; t \in d\}|}] \qquad (7)$$

$$tf - idf(t, d, D) = tf(t, d).idf(t, D) \qquad (8)$$

Where $f_d(t)$ is the rate of the term 't' in the document 'd', and 'w' is a set of words in the document 'd' while 'D' is the corpus of documents. The TF-IDF model combines Term Frequency (TF) with Inverse Document Frequency (IDF) to measure the importance of a term in a document relative to a corpus. IDF indicates how common or rare a term is across all documents and is calculated using a logarithmic formula, as shown in Eq. 7. The final TF-IDF model is presented in Eq. 8.

In Text Classification (TC), representation and feature extraction processes often lead to many terms, causing issues like the "*curse of dimensionality*" This is due to the inherent characteristics of textual data, like noise, redundancy, and sparseness. To address this, selecting an efficient Feature Selection (FS) technique is crucial. The Chi-Square (C.H.I.) method is commonly used to choose relevant features in various works.

### E. Arabic Text Classification

TC is the task of deciding whether a piece of text belongs to prescribed classes based on understanding and discriminating the pattern of the text [29]. It is a significant task involved in the TC system [3] [30]. TC is one of the most challenging computational tasks in the ML community. Only a few researchers worked in ATC, and numerous classifier learning algorithms have been used, such as Naive Bayes (NB) [31], Support Vector Machine (SVM) [4] and Artificial Neural Network (ANN) [32].

*1) Bag of Words (BoW):* BoW is a textual presentation style suitable for classification models in which the text is treated as a collection of words, regardless of syntax. BoW shows whether or not a certain word occurs in the document, with no regard for word order. The performance of the various ML algorithms we developed using BoW was unsatisfactory owing to the loss of semantic and syntactic information between phrases. To boost performance, we used different representation schemes that can accept semantic and syntactic differences.

*2) Term Frequency—Inverse Document Frequency (TFIDF):* Term Frequency (TF) is a popular textual presentation approach that is equivalent to the BoW technique. The term's recurrence in a given text is what determines TF, whereas its existence determines BoW. Eq. 9 represents the frequency of any word in the supplied document. However, typical terms included in all documents, such as articles, conjunctions, and prepositions, receive a low rating since they add little to the content. Thus, we employ Inverse Document Frequency (IDF) to reduce the value of terms that appear frequently in the document collection while increasing the significance of phrases that appear infrequently. IDF is consistent across corpora and determines the proportion of papers that have a certain phrase. Eq. 10 describes how it is assessed. TF-IDF is a statistical criterion for determining how closely a phrase is associated with a document in a collection of manuscripts, known as a corpus. TF-IDF is calculated by multiplying TF and IDF. TFIDF is a straightforward technique to text categorization. Thus, the TF-IDF is created during model training and subsequently applied to the test set.

$$TF = \frac{Number \ of \ times \ a \ term \ appear \ in \ the \ document}{Total \ number \ of \ terms \ in \ the \ document} \qquad (9)$$

$$IDF = Log_{10} \frac{Total \ number \ of \ Documents}{Number \ of \ documents \ that \ includes \ the \ term} \qquad (10)$$

### F. Building of Classification Models (Learning)

*1) Multinomial Naive Bayes (MNB):* The naive Bayes classifier is based on Bayes Theorem, which operates on conditional probability. The conditional probability is the likelihood that something will happen if something else has already happened [15]. Eq. 11 provides the formula for computing conditional probability, with *A* representing the hypothesis and *B* representing the evidence.

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)} \qquad (11)$$

NB computes numerous model parameters for a given set of labeled training data, including the likelihood of each class label happening. Then, estimate the class for each set of test data based on its chance of being assigned to different classes. The expected class is determined by the largest likelihood [33]. Multinomial Naive Bayes (MNB) is a prominent supervised learning classification approach used to analyze categorical text data. It is a probabilistic learning strategy that is widely used

in natural language processing (NLP). To anticipate a text's tag, the algorithm applies the Bayes principle. It calculates the likelihood of each tag for a given sample and presents the tag with the highest probability as output.

*2) Bernoulli Naive Bayes (BNB):* It is a subset of the Naive Bayes Algorithm that is used to classify binary features such as '1' or '0' that are independent of one another. Bernoulli Naive Bayes is used to identify spam, classify text, do sentiment analysis, and determine whether a given word occurs in a document. Eq. 12 expresses the decision rule of Bernoulli NB, where $P(x_i|y)$ is the conditional probability of $x_i$ happening if $y$ has occurred, $i$ is the event, and $x_i$ has a binary value of 0 or 1.

$$P(x_i|y) = P(x_i = 1|y)x_i + (1 - P(x_i = 1|y))(1 - x_i) \quad (12)$$

*3) Stochastic Gradient Descent (SGD):* Gradient descent is an optimization approach for training machine learning models and neural networks. The training data allows these models to learn over time, and the cost function in gradient descent quantifies accuracy with each parameter update, i.e., direction and learning rate. The model will keep modifying its parameters to get the minimal feasible error. Stochastic gradient descent (SGD) handles one training epoch per example and changes its parameters one at a time. SGDs are easy to retain in memory since they only require one training sample. Its regular updates assist in avoiding the local minimum and discovering the global one.

*4) Support Vector Classifier (SVC):* It is a specific implementation of the Support Vector Machine (SVM) algorithm developed for classification. It aims to discover the n-dimensional hyperplanes that sufficiently divide the data instances into various types.

*5) Logistic Regression (LR):* Regression modeling is a well-known and reliable statistical approach for analyzing and describing the relationship between a dependent variable and a group of independent predictors. Logistic regression is a specific case in regression modeling in which the result is binary. A logistic function, also known as a *sigmoid function*, may be used to explain logistic regression. This function receives any actual input x and returns a probability value between 0 and 1 as defined in Eq. 13.

$$f(x) = \frac{1}{1 + e^{-x}} \quad (13)$$

*6) Linear Support Vector Classifier (LSVC):* The Linear Support Vector Classifier (SVC) approach uses a linear kernel function to complete classification and it is more suitable than SVC for large datasets.

## V. EXPERIMENTAL RESULTS AND COMPARATIVE ANALYSIS

The machine learning algorithms have been implemented using Python version 3.8.0 in the Anaconda environment, specifically within the Jupyter Notebook. Python machine-learning libraries such as NLTK, pandas, and sci-kit-learn are utilized to evaluate the proposed methods' performance. The findings and discussions related to the various methods employed are delineated in the following sections. We used the

sci-kit-learn library, which contains ML algorithms for the experimentation. To evaluate the effect of the preprocessing and representation on classification, we conducted several experiments on the BBC Arabic dataset and COVID-19 and analyzed the performance results of various classifiers. Additionally, the proposed is evaluated using classifications generated from the baseline model. We summarise all experimental results and compare the proposed method with other methods.

In general, the works provide a comprehensive summary of the performance results for two datasets with preprocessing and without feature selection and vice versa, allowing for comparing different transformation methods based on various evaluation metrics. We found that pre-processing, representation of AT, and feature engineering techniques played an essential role in enhancing the performance of ATC. In the beginning, pre-processing can affect ATC's performance. In addition, there were other techniques of representation, such as BoW and TF-IDF, which have some drawbacks, like missing the order of the words and losing the meaning of the words. Pre-processing techniques, such as stop word removal, can be used to reduce the dimension, but in some cases, pre-processing can positively or negatively affect the performance. Finally, we know accuracy is insufficient to evaluate ATCs, so we extend our investigation using different evaluation metrics. We summarise our findings in the conclusion section after the results and discussions.

### A. Results on COVID-19 Dataset

This section discusses the experimental result of the COVID-19 dataset. It evaluates the model's performance and studies its effects on ATC regarding Accuracy (ACC), Precision (PR), Recall (RE), and F1-score. We discuss and evaluate the experimental result based on different methods.

First, we applied six classification algorithms on the COVID-19 dataset for ATC without pre-processing and without feature selection. Table III shows the evaluation metrics for this scenario, where the accuracy ranges from 81% to 83%. Regarding the macro metrics, the precision ranges from 52% to 69% while the recall ranges from 34% to 53%. F1-score ranges from 31% to 57%. Generally, the weighted metrics show a better performance since the contribution of each class is considered where the precision ranges from 75% to 80% while the recall ranges from 81% to 83%. Weighted F1-score ranges from 73% to 80%. Generally, SGDC and LSVC have the best performance among the used classifiers.

TABLE III. EVALUATION METRICS FOR COVID-19 DATASET **WITHOUT** PREPROCESSING AND **WITHOUT** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | 81 | 52 | 34 | 31 | 75 | 81 | 73 |
| BNBC | 81 | 52 | 34 | 31 | 75 | 81 | 73 |
| LRC | 82 | **69** | 43 | 47 | 79 | 82 | 78 |
| SGDC | **83** | 68 | 51 | 55 | **80** | **83** | **80** |
| SVC | 82 | 67 | 42 | 45 | 78 | 82 | 77 |
| LSVC | 82 | 66 | **53** | **57** | **80** | 82 | **80** |

Table IV shows the evaluation metrics for classification algorithms on the COVID-19 dataset with pre-processing but without feature selection. The LSVC has the highest accuracy of 83%. The NBC and BNBC have the lowest accuracy of

81%. Also, they have the lowest metrics with a precision of 52%, recall of 34%, and F1-Score of 31% for macro metrics. Their weighted metrics are minimal, i.e., precision, recall, and F1-Score are 75%, 81%, and 73%, respectively. The LSCV almost achieves the best performance indicated by both macro and weighted metrics.

TABLE IV. EVALUATION METRICS FOR COVID-19 DATASET **WITH** PREPROCESSING AND **WITHOUT** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | 81 | 52 | 34 | 31 | 75 | 81 | 73 |
| BNBC | 81 | 52 | 34 | 31 | 75 | 81 | 73 |
| LRC | 82 | **70** | 43 | 47 | 79 | 82 | 77 |
| SGDC | 82 | 68 | 50 | 55 | **80** | 82 | 80 |
| SVC | 82 | 67 | 42 | 45 | 78 | 82 | 77 |
| LSVC | **83** | 67 | **52** | **57** | **80** | **83** | **81** |

As shown in Table V, five classifiers have an 82% accuracy when the COVID-19 dataset is evaluated without preprocessing but with feature selection. Regarding the macro metrics, the precision ranges from 59% to 77% while the recall ranges from 23% to 52%. F1-score ranges from 35% to 54%. For the weighted metrics, the precision ranges from 77% to 81% while the recall is either 78% or 82%. Weighted F1-score ranges from 74% to 78%. There is no single classifier that is the best in most metrics, where the best classifiers are BNBC, SGDC, and LSVC. As with the previous two scenarios, the weighted metrics are higher than the macro ones because the contribution/weight of each class is considered.

TABLE V. EVALUATION METRICS FOR COVID-19 DATASET **WITHOUT** PREPROCESSING AND **WITH** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | **82** | 74 | 23 | 35 | 79 | **82** | 74 |
| BNBC | 78 | 59 | **52** | 54 | 79 | 78 | **78** |
| LRC | **82** | 71 | 42 | 44 | 79 | **82** | 77 |
| SGDC | **82** | **77** | 41 | 43 | **81** | **82** | 76 |
| SVC | **82** | 60 | 45 | 48 | 77 | **82** | 77 |
| LSVC | **82** | 66 | 46 | 51 | 78 | **82** | **78** |

The last scenario is when the COVID-19 dataset is classified and the metrics evaluated with data preprocessing and with feature selection as given in Table VI. The accuracy ranges from 79% to 82%. Regarding the macro metrics, the precision ranges from 55% to 81% while the recall ranges between 37% and 52%. The minimum F1-score is 36% and the maximum is 53%. When the metrics are evaluated with the weight of classes, the precision ranges between 77% and 81% while the recall is either 78% or 82%. The minimum F1-score is 74% and the maximum is 78%. The best classifiers are BNBC, SGDC, and LSVC.

From the Tables III, IV, V, and VI we notice that the weighted metrics are always better than the macro metrics. Also, there is no specific classifier that is the best in the seven metrics, i.e., accuracy, macro and weighted (precision, recall, and F1-score). In most of the cases, SGDC and LSVC were the best.

### B. Results on BBC Dataset

We extended our experimentations in this section by studying the performance of long ATC using the BBC dataset. The

TABLE VI. EVALUATION METRICS FOR COVID-19 DATASET **WITH** PREPROCESSING AND **WITH** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | **82** | 81 | 37 | 36 | 81 | 82 | 74 |
| BNBC | 79 | 55 | **52** | **53** | 77 | 79 | **78** |
| LRC | **82** | 78 | 39 | 41 | 80 | **82** | 76 |
| SGDC | **82** | 78 | 38 | 39 | **81** | **82** | 75 |
| SVC | 81 | 63 | 40 | 42 | 77 | 81 | 76 |
| LSVC | **82** | 68 | 42 | 45 | 79 | **82** | 77 |

study evaluates six classification algorithms on an extended BBC Arabic dataset for Arabic text classification (ATC).

Table VII shows the evaluation metrics for the classification of the BBC dataset for ATC without pre-processing and without feature selection. The accuracy ranges from 68% to 93%, where the accuracy of SGDC is 93% and the accuracy of LSVC is 92%. Regarding the macro metrics, the precision ranges from 78% to 94% where SVC has the maximum value and both SGDC and LSVC have a value of 93%. The maximum recall is 89% while the minimum is 34%. The F1-score ranges from 38% and 91% where the SGDC has the maximum value while LSVC has a score of 90%. For the weighted metrics, the SGDC has the best precision, recall, and F1-score with a value of 93% for all of them, while LSVC has a value of 92% for the three metrics. Thus, the best classifier is the SGDC followed by LSVC while both NBC and BNBC show the minimum metrics.

TABLE VII. EVALUATION METRICS FOR BBC DATASET **WITHOUT** PREPROCESSING AND **WITHOUT** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | 68 | 78 | 34 | 38 | 76 | 68 | 63 |
| BNBC | 68 | 78 | 34 | 38 | 76 | 68 | 63 |
| LRC | 86 | 80 | 62 | 68 | 86 | 86 | 85 |
| SGDC | 93 | 93 | 89 | 91 | 93 | 93 | 93 |
| SVC | 87 | 94 | 68 | 76 | 88 | 87 | 86 |
| LSVC | 92 | 93 | 88 | 90 | 92 | 92 | 92 |

The metrics for ATC for the BBC dataset with preprocessing and without feature selection are shown in Table VIII. There is a great similarity with the values reported in Table VII. The best classifier is the SGDC followed by LSVC, while both NBC and BNBC show the minimum metrics.

TABLE VIII. EVALUATION METRICS FOR BBC DATASET **WITH** PREPROCESSING AND **WITHOUT** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | 69 | 78 | 35 | 39 | 77 | 69 | 64 |
| BNBC | 69 | 78 | 35 | 39 | 77 | 69 | 64 |
| LRC | 86 | 80 | 61 | 67 | 86 | 86 | 85 |
| SGDC | 93 | 94 | 89 | 91 | 93 | 93 | 93 |
| SVC | 87 | 94 | 69 | 77 | 88 | 87 | 86 |
| LSVC | 92 | 93 | 88 | 90 | 92 | 92 | 92 |

The results of classification without preprocessing but with feature selection are shown in Table IX. NBC has a minimal accuracy of 57% while the SVC has a maximum accuracy of 97%. Both BNBC and SGDC have an accuracy of 96% while the LSVC has a 95% accuracy. For the metrics that are evaluated at the macro level, the precision ranges from 50%

to 96%, where both SGDC and LSVC have a value of 96% and the SVC has a precision of 95%. The macro recall ranges from 23% to 96%, while the F1-score is between 24% and 95%. The weighted precision ranges from 66% to 98%. The SVC has the maximum precision while both BNBC and SGDC have a precision of 96%. The weighted recall ranges from 57% and 97%. Both BNBC and SGDC have a precision of 96% and NBC has the minimal value. The range of the F1-score is similar to the precision where the SVC has the maximum of 97% followed by 96% for both BNBC and SGDC, while NBC has the lowest F1-score. The worst classifier regarding all metrics is NBC while the best classifier is the SVC, followed by SGDC and LSVC.

TABLE IX. EVALUATION METRICS FOR BBC DATASET **WITHOUT** PREPROCESSING AND **WITH** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | 57 | 50 | 23 | 24 | 66 | 57 | 47 |
| BNBC | 96 | 94 | 88 | 88 | 96 | 96 | 96 |
| LRC | 71 | 78 | 42 | 48 | 78 | 71 | 67 |
| SGDC | 96 | **96** | 91 | 93 | 96 | 96 | 96 |
| SVC | **97** | 95 | **96** | **95** | 98 | 97 | **97** |
| LSVC | 95 | **96** | 89 | 92 | 95 | 95 | 95 |

The metrics for BBC dataset classification with preprocessing and with feature selection are shown in Table X. Similar to Table IX, the best classifier is SVC followed by SGDC, LSVC, and BNBC. The minimal metrics are obtained by NBC.

We notice that without feature selection, as shown in Table VII and VIII, the best classifier is SGDC followed by LSVC. However, with feature selection, as shown in Table IX and X, the best classifier is SVC followed by SGDC.

TABLE X. EVALUATION METRICS FOR BBC DATASET **WITH** PREPROCESSING AND **WITH** FEATURE SELECTION

| Transformation Method | Acc | Macro | | | Weighted | | |
|---|---|---|---|---|---|---|---|
| | | PR | RE | F1-score | PR | RE | F1-score |
| NBC | 49 | 07 | 14 | 09 | 24 | 49 | 32 |
| BNBC | 95 | 92 | 87 | 87 | 96 | 95 | 95 |
| LRC | 78 | 38 | 34 | 35 | 71 | 78 | 72 |
| SGDC | 96 | **97** | 89 | 92 | 96 | 96 | 96 |
| SVC | **98** | 96 | **96** | **96** | 99 | 98 | **98** |
| LSVC | 94 | **97** | 84 | 88 | 95 | 94 | 94 |

## VI. DISCUSSION AND RESULTS

This section discusses the effectiveness of preprocessing and feature selection on short and long Arabic Text classification. The authors of [34] evaluated representation and preprocessing on a short text dataset and indicated that the effectiveness of preprocessing is positive in the case with Bow and negative in others with TFID. Still, LRC and SVC generally offered the best performance across most metrics. The authors found that as the number of characteristics increased, so did the execution time. In the meanwhile, the classifiers' performance remains unchanged. After increasing the number of characteristics from 7,000 to 10,000, all classifiers except SVC maintained their accuracy.

The experimentation evaluates the performance of various models on short data for ATC without any preprocessing or feature selection. We observed under these conditions that

all models achieved accuracy scores between 81% and 83%. Among them, LRC had the best macro accuracy, while LSVC had the highest macro recall and F1-score. When weighted criteria like precision, recall, and F1-Score were considered, SGDC and LSVC emerged as standouts.

When introducing preprocessing, but still without feature selection, LSVC consistently outperformed the other models across almost all metrics. On the other hand, NBC and BNBC demonstrated the lowest performance.

The experimental outcomes underscore that preprocessing had mixed effects on model results while it improved performance for some models, it hindered the performance of others. LSVC and SGDC were the top-performing models across various scenarios, while NBC and BNBC lagged. When we continued our experimentation on the BBC dataset, in the initial experiment, which focused on the dataset without preprocessing and feature selection, the SGDc and LSVC with TF-IDF were the top performers. Specifically, SGDc and LSVC achieved the highest accuracy, with 93% and 92%, respectively. In terms of macro-precision, both algorithms scored 93%. For macro-recall and macro-F-score, the top scores were 89% and 91%, respectively. In weighted metrics, SGDC outperformed with 93% across all metrics, followed closely by LSVC at 92%.

The subsequent experiment, which incorporated preprocessing but left out feature selection, affirmed the superiority of the SGDC and LSVC classifiers for the extended BBC dataset. The third experiment did not involve preprocessing while including feature selection. The SVC stood out by almost securing the highest values. Notably, the NBC and LRS underperformed compared to the others in this context. The final experiment, including preprocessing and feature selection, yielded comparable results to the third experiment. The variations among metrics were minimal, between 1% and 3%. SVC was a standout again. Many models performed exceptionally well in weighted metrics, such as SGDC and LSVC. However, as in the third experiment, NBC and LRC lagged the rest.

In summary, when working with the BBC dataset for Arabic text classification, the SGDC and LSVC consistently demonstrate high performance, especially without preprocessing and feature selection. However, with preprocessing and feature selection, SVC tends to be the best performer, while NBC and LRC trail behind their counterparts. The experiment's result demonstrated that the preprocessing and the feature selection impact the text classification performance, and the dataset type (short/long) also weighs heavily on the performance.

## VII. CONCLUSION

In categorizing Arabic text, a complete study investigation was undertaken to demonstrate the usefulness of pre-processing, feature extraction, feature selection, and the dataset's features. Numerous ML models have been introduced to underscore the efficacy of diverse techniques in classifying Arabic text. The study's results indicate that many models influence the accuracy of system performance in ATC. The experimentation indicates that representation, encompassing feature extraction and feature selection, is essential in ATC.

Simultaneously, the preprocessing, classification algorithm and the dataset's characteristics influence the classification performance's efficacy. The findings clearly illustrate the benefits of the feature representation method and its impact on text classification efficacy. Based on the analysis presented in this article, which is limited to two datasets, several outstanding issues remain for future research, such as the absence of benchmark datasets, the shortage of lexicons, and the challenge of identifying techniques that address the contextual significance of ATC. The study highlights the effectiveness of pre-processing and feature representation on text classification performance. Challenges remain, like a lack of benchmark datasets and context-aware techniques for ATC; opportunities exist for enhancing tools like data augmentation and preprocessing techniques, mainly stemming.

REFERENCES

[1] M. Allahyari, S. Pouriyeh, M. Assefi, S. Safaei, E. D. Trippe, J. B. Gutierrez, and K. Kochut, "A brief survey of text mining: Classification, clustering and extraction techniques," *arXiv preprint arXiv:1707.02919*, 2017.

[2] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[3] F. Sebastiani, "Machine learning in automated text categorization," *ACM computing surveys (CSUR)*, vol. 34, no. 1, pp. 1–47, 2002.

[4] A. Moh'd A Mesleh, "Chi square feature extraction based svms arabic language text categorization system," *Journal of Computer Science*, vol. 3, no. 6, pp. 430–435, 2007.

[5] H. Alsayadi, A. Abdelhamid, I. Hegazy, and Z. Taha, "Data augmentation for arabic speech recognition based on end-to-end deep learning," *International Journal of Intelligent Computing and Information Sciences*, vol. 21, no. 2, pp. 50–64, 2021.

[6] F. Sadat, F. Kazemi, and A. Farzindar, "Automatic identification of arabic dialects in social media," in *Proceedings of the first international workshop on Social media retrieval and analysis*, 2014, pp. 35–40.

[7] N. Y. Habash, *Introduction to Arabic natural language processing.* Springer Nature, 2022.

[8] H. Tavasoli, B. J. Oommen, and A. Yazidi, "On utilizing weak estimators to achieve the online classification of data streams," *Engineering Applications of Artificial Intelligence*, vol. 86, pp. 11–31, 2019.

[9] M. Suhil, "Representation and classification of text data," 2019.

[10] A. Ayedh, G. Tan, K. Alwesabi, and H. Rajeh, "The effect of preprocessing on arabic document categorization," *Algorithms*, vol. 9, no. 2, p. 27, 2016.

[11] I. Guellil, H. Saâdane, F. Azouaou, B. Gueni, and D. Nouvel, "Arabic natural language processing: An overview," *Journal of King Saud University-Computer and Information Sciences*, vol. 33, no. 5, pp. 497–507, 2021.

[12] S. L. Marie-Sainte, N. Alalyani, S. Alotaibi, S. Ghouzali, and I. Abunadi, "Arabic natural language processing and machine learning-based systems," *IEEE Access*, vol. 7, pp. 7011–7020, 2018.

[13] M. Masadeh, A. Masadeh, O. Alshorman, F. Khasawneh, and M. Masadeh, "An efficient machine learning-based covid-19 identification utilizing chest x-ray images," *IAES International Journal of Artificial Intelligence*, pp. 356–366, 2022.

[14] M. Masadeh, O. Hasan, and S. Tahar, "Machine-learning-based self-tunable design of approximate computing," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 29, no. 4, pp. 800–813, 2021.

[15] M. Masadeh, A. Aoun, O. Hasan, and S. Tahar, "Decision tree-based adaptive approximate accelerators for enhanced quality," in *International Systems Conference (SysCon)*. IEEE, 2020, pp. 1–5.

[16] M. Masadeh, O. Hasan, and S. Tahar, "Machine learning-based self-compensating approximate computing," in *2020 IEEE International Systems Conference (SysCon)*. IEEE, pp. 1–6.

[17] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning.* MIT press, 2016.

[18] R. NL, "Stopword lists," 2024, https://www.ranks.nl/stopwords/arabic [Accessed: 17.01.2024].

[19] A. K. Uysal and S. Gunal, "The impact of preprocessing on text classification," *Information processing & management*, vol. 50, no. 1, pp. 104–112, 2014.

[20] K. Kowsari, K. Jafari Meimandi, M. Heidarysafa, S. Mendu, L. Barnes, and D. Brown, "Text classification algorithms: A survey," *Information*, vol. 10, no. 4, p. 150, 2019.

[21] S. Minaee, N. Kalchbrenner, E. Cambria, N. Nikzad, M. Chenaghlu, and J. Gao, "Deep learning–based text classification: a comprehensive review," *ACM computing surveys (CSUR)*, vol. 54, no. 3, pp. 1–40, 2021.

[22] M. Masadeh, H. J. Davanager, and A. Y. Muaad, "A novel machine learning-based framework for detecting religious arabic hatred speech in social networks," *International Journal of Advanced Computer Science and Applications*, vol. 13, no. 9, 2022.

[23] A. Y. Muaad, H. Jayappa, M. A. Al-antari, and S. Lee, "Arcar: a novel deep learning computer-aided recognition for character-level arabic text representation and recognition," *Algorithms*, vol. 14, no. 7, p. 216, 2021.

[24] Y. Albalawi, J. Buckley, and N. S. Nikolov, "Investigating the impact of pre-processing techniques and pre-trained word embeddings in detecting arabic health information on social media," *Journal of big Data*, vol. 8, no. 1, p. 95, 2021.

[25] Y. A. Alhaj, J. Xiang, D. Zhao, M. A. Al-Qaness, M. Abd Elaziz, and A. Dahou, "A study of the effects of stemming strategies on arabic document classification," *IEEE access*, vol. 7, pp. 32 664–32 671, 2019.

[26] M. K. Saad and W. Ashour, "Osac: Open source arabic corpora," in *6th ArchEng Int. Symposiums, EEECS*, vol. 10, 2010, p. 55.

[27] H. Benjamin and S. Sotardi, "The pareto principle," *J Am College Radiol*, vol. 15, no. 6, p. 931, 2018.

[28] T. Kanan, B. Hawashin, S. Alzubi, E. Almaita, A. Alkhatib, K. A. Maria, and M. Elbes, "Improving arabic text classification using p-stemmer," *Recent Advances in Computer Science and Communications (Formerly: Recent Patents on Computer Science)*, vol. 15, no. 3, pp. 404–411, 2022.

[29] G. Kanaan, R. Al-Shalabi, S. Ghwanmeh, and H. Al-Ma'adeed, "A comparison of text-classification techniques applied to arabic text," *Journal of the American society for information science and Technology*, vol. 60, no. 9, pp. 1836–1844, 2009.

[30] A. Hotho, A. Nürnberger, and G. Paaß, "A brief survey of text mining," *Journal for Language Technology and Computational Linguistics*, vol. 20, no. 1, pp. 19–62, 2005.

[31] S. Ganjavi, P. Georgiou, and S. Narayanan, "A transcription scheme for languages employing the arabic script motivated by speech processing applications," in *Proceedings of the Workshop on Computational Approaches to Arabic Script-based Languages*, 2004, pp. 59–65.

[32] F. Harrag and E. El-Qawasmah, "Neural network for arabic text classification," in *Second International Conference on the Applications of Digital Information and Web Technologies*. IEEE, 2009, pp. 778–783.

[33] M. K. A. Aljero and N. Dimililer, "A Novel Stacked Ensemble for Hate Speech Recognition," *Applied Sciences*, vol. 11, no. 24, p. 11684, 2021.

[34] A. Y. Muaad, H. J. Davanagere, D. Guru, J. B. Benifa, C. Chola, H. AlSalman, A. H. Gumaei, and M. A. Al-antari, "Arabic document classification: performance investigation of preprocessing and representation techniques," *Mathematical Problems in Engineering*, vol. 2022, pp. 1–16, 2022.

# Evaluating Tree-based Ensemble Strategies for Imbalanced Network Attack Classification

Hui Fern Soon[1], Amiza Amir[2], Hiromitsu Nishizaki[3], Nik Adilah Hanin Zahri[4],
Latifah Munirah Kamarudin[5], Saidatul Norlyana Azemi[6]
Faculty of Electronic Engineering & Technology
Universiti Malaysia Perlis Arau, Perlis, Malaysia
University of Yamanashi, Kofu, Yamanashi, Japan[1]
Faculty of Electronic Engineering & Technology
Centre of Excellence for Advanced Computing (ADVCOMP)
Universiti Malaysia Perlis Arau, 02600, Perlis, Malaysia[2,4]
Integrated Graduate School of Medicine, Engineering & Agricultural Science
University of Yamanashi, Kofu, Yamanashi, Japan[3]
Faculty of Electronic Engineering & Technology
Centre of Excellence for Advanced Sensor Technology (CEASTECH)
Universiti Malaysia Perlis Arau, 02600, Perlis, Malaysia[5]
Faculty of Electronic Engineering & Technology
Centre of Excellence for Advanced Communication Engineering (ACE)
Universiti Malaysia Perlis Arau, 02600, Perlis, Malaysia[6]

*Abstract*—With the continual evolution of cybersecurity threats, the development of effective intrusion detection systems is increasingly crucial and challenging. This study tackles these challenges by exploring imbalanced multiclass classification, a common situation in network intrusion datasets mirroring real-world scenarios. The paper aims to empirically assess the performance of diverse classification algorithms in managing imbalanced class distributions. Experiments were conducted using the UNSW-NB15 network intrusion detection benchmark dataset, comprising ten highly imbalanced classes. The evaluation includes basic, traditional algorithms like the Decision Tree, K-Nearest Neighbor, and Gaussian Naive Bayes, as well as advanced ensemble methods such as Gradient Boosted Decision Trees (GraBoost) and AdaBoost. Our findings reveal that the Decision Tree surpassed the Multi-Layer Perceptron, K-Nearest Neighbor, and Naive Bayes in terms of overall F1-score. Furthermore, thorough evaluations of nine tree-based ensemble algorithms were performed, showcasing their varying efficacy. Bagging, Random Forest, ExtraTrees, and XGBoost achieved the highest F1-scores. However, in individual class analysis, XGBoost demonstrated exceptional performance relative to the other algorithms. This is confirmed by achieving the highest F1-scores in eight out of the ten classes within the dataset. These results establish XGBoost as a predominant method for handling multiclass imbalance classification with Bagging being the closest feasible alternative, as Bagging gains an almost similar accuracy and F1-score as XGBoost.

*Keywords—Multiclass imbalanced classification; ensemble algorithm; network attack; UNSW-NB15 dataset; F1-score*

## I. INTRODUCTION

Following the COVID-19 pandemic, accelerated advancements in information technology have reshaped organizational operations, interpersonal interactions, and service delivery methods. The Internet and cyber technology have facilitated a highly interconnected global society, significantly influencing almost every facet of the modern world. This revolutionizes human lifestyles, transforms various industries, and promotes global innovation. The shift towards remote work and virtual platforms has surged, prompting the development of new tools and technologies to accommodate these changes. Additionally, the healthcare sector has experienced a growth in telemedicine and digital health solutions, enabling remote patient consultations and monitoring. However, these advancements also increase vulnerability to cybersecurity attacks, as cybercriminals view the rapid expansion of IT applications, especially in e-commerce, as lucrative targets. The European Union Agency for Cybersecurity (ENISA) noted a notable rise in cybersecurity incidents during the latter part of 2022 and the first half of 2023, as referenced in [1]. These developments underscore the urgent need for effective, reliable, and robust defense systems against such attacks. Concurrently, with the proliferation of AI, machine learning and deep learning algorithms have emerged as powerful tools for network security.

The effectiveness of machine learning and deep learning models in detecting network attacks hinges on the quality and relevance of the training data. Inadequate or irrelevant training data can yield inaccurate or unreliable outcomes. Therefore, it is essential to ensure the training data for these models is high-quality and representative of actual network attack scenarios. Typically, network traffic remains normal until a cyberattack or network failure occurs, causing a deviation from usual patterns. Machine learning and deep learning models are capable of identifying and learning these anomalies, thereby precisely detecting and classifying network attacks.

Consequently, most of the training data will consist of normal network traffic. The abnormal network traffic dataset, representing potential network attacks, includes various categories of network assaults. Rare or novel attack types might have limited sample sizes, potentially smaller than those found

in common attack types or normal traffic data. This leads to a significant imbalance in class composition, potentially introducing model bias, rendering predictions unreliable, and hindering the detection of rare or new attacks. As noted by [2], most network intrusion datasets are inherently multiclass imbalanced, reflecting real-world conditions. In such datasets, class distribution is uneven (as depicted in Fig. 1), with some classes being minority and others majority.



Fig. 1. Example of multiclass imbalanced classification.

Imbalanced datasets can skew classifiers, biasing them toward the majority class [3]. This presents a significant challenge, as most classifiers are inherently designed for balanced scenarios. One simplistic approach could be to exclude minority classes with insufficiently sized samples from the dataset. However, this could result in models that are outdated with respect to the latest cyberattacks. We continuously update and retrain these models with fresh data to enhance their adaptability to evolving attack techniques and ensure sustained effectiveness.

Addressing the imbalance often involves sampling solutions. Techniques such as Random Oversampling [4] or Synthetic Minority Oversampling Technique (SMOTE) [5] augment infrequent cases, while methods like Random Undersampling [6] or Tomek links [7] reduce redundancies in the dataset by decreasing majority samples. Hybrid techniques like SMOTEENN [8], which combine oversampling and undersampling, and ROSE (Random OverSampling Examples) [9], which create synthetic spaces between classes, are also utilized. However, oversampling risks overfitting, and undersampling may lead to information loss. SMOTEENN and ROSE, while versatile, are also prone to overfitting. Moreover, the continually changing nature of new attacks complicates the use of these methods, given the dynamic class distributions. Thus, these methods can temporarily achieve balance but have limitations in long-term applicability and robustness. Consequently, this paper does not focus on sampling solutions but rather on the inherent capabilities of classification algorithms to address imbalanced class problems effectively.

While various machine learning approaches have been proposed for network attack classification, a predominant focus remains on enhancing overall accuracy—a metric poorly suited for imbalanced multiclass datasets. Accuracy measures the proportion of correct predictions made by the model, but it fails to adequately represent minority classes, particularly those with low sample sizes. An accuracy-centric model might disregard minority classes, classifying all instances as the majority class, thereby achieving high overall accuracy but poor detection of rare, yet critical, cases. This oversight necessitates

a more nuanced, class-specific evaluation. Additionally, there is a notable research gap concerning the effectiveness of different algorithms in addressing multiclass imbalances.

Therefore, this research has a dual focus. First, it seeks to identify which conventional machine learning algorithms are best suited for addressing the unique challenges of multiclass imbalanced classification, specifically in the context of network attack classification. Second, it explores which ensemble algorithms are most effective in these scenarios. Following guidance from [10], potential solutions include sampling techniques, ensemble methods, cost-sensitive learning, and deep learning methods. This paper, however, concentrates on the application of ensemble approaches to manage imbalanced data scenarios. We compare and experiment with a range of machine learning algorithms, from simpler ones like decision trees and K-nearest neighbors to more complex ensemble algorithms such as Gradient Boosted Decision Trees (GraBoost) and AdaBoost. The objective is to ascertain the most effective algorithm for addressing the complexities of imbalanced datasets in network intrusion detection.

The experimental evaluation utilizes the publicly available UNSW-NB15 dataset [11], characterized by a highly imbalanced class distribution. Initial experiments compared the performance of a single Decision Tree (DT) against instance-based methods like K-Nearest Neighbor (KNN), function-based models including Multilayer Perceptron (MLP), and Bayesian-based approaches exemplified by Naive Bayes (NB).

Despite the initial success of the Decision Tree, there is a need for more precision, particularly in identifying tree-based ensemble algorithms that excel in multiclass imbalance classification. This research thus focuses on discovering the most effective tree-based ensemble algorithms for managing the challenges posed by imbalanced multiclass datasets. In addition to a single Decision Tree, we conducted experiments comparing nine tree-based ensemble learning algorithms: Bagging with a Decision Tree as the base classifier, Random Forest (RF), Extremely Randomized Trees (ExtraTree), Adaptive Boosting (AdaBoost), Gradient Boosting (GraBoost), Histogram-based Gradient Boosting (HistGraBoost), Extreme Gradient Boosting (XGBoost), Light Gradient Boosting Machine (LightGBM), and Categorical Gradient Boosting (CatBoost).

To summarize, the paper's primary contributions are as follows. First, preliminary results indicate the superiority of the Decision Tree over other traditional machine learning algorithms. Second, XGBoost has been determined as the optimal tree-based ensemble method for multi-class imbalanced classification with Bagging being the closest feasible alternative. Third, this paper offers practitioners a powerful approach to address the issues often encountered with imbalanced multi-class datasets effectively. Consequently, this improves the overall efficacy of cybersecurity protocols.

The structure of this paper is as follows. Section II explains the related work. Section II describes the methodology. Section IV illustrates the dataset, algorithms and performance metrics used in this research, while Section V describes results of the algorithms. Finally, in Section VI the conclusions and the future works are being discuss.

## II. LITERATURE REVIEW

Network attack detection datasets are often multiclass imbalanced [2]. Nevertheless, despite this observed pattern, many research efforts continue to pay attention to tackling the issue of imbalanced classification problems. Typical solutions to dealing with imbalanced dataset issues include utilising sampling approaches [12], [13]. The first sampling approach is to apply the oversampling technique to address the imbalance in minority classes. Random Over-sampling involves the random duplication of cases from the minority class [4]. SMOTE [5], which stands for Synthetic Minority Over-sampling Technique, is a method used to create synthetic samples comparable to the minority data cluster. The second sampling approach is by under-sampling majority classes. For example, the Random Under-sampling [6] randomly removes the majority of class examples, and Tomek links [7] work by removing overlap between class sample distributions. Finally, hybrid/ensemble sampling refers to a technique that combines multiple sampling methods or models to improve the accuracy and reliability of the sampling process. For instance, SMOTEENN [8] is a technique that combines SMOTE over-sampling with edited closest neighbour under-sampling. ROSE(Random OverSampling Examples) sampling [9] generates smooth distributions by creating synthetic spaces between minority and majority examples.

Nevertheless, the deliberate process of oversampling minority classes can lead to over fitting of the model due to the replication and noise. On the contrary, by undersampling the majority classes, there is a risk of losing valuable information crucial for precise classification. Hybrid or ensemble sampling techniques, such as SMOTEEN and Rose sampling, prove helpful in generating more balanced sample distributions. However, class distribution patterns in complex real-world contexts are rarely uniform or evenly distributed. In addition, they inherit the overfitting and losing valuable issues from oversampling and undersampling, respectively. Another key challenge is class distribution concept drift in the dynamic network traffic data. It is possible that novel network attacks will emerge, each with a limited sample size. As relative class frequencies change over time, a previously balanced data set may become outdated.

Hence, this paper aims to identify the best algorithm without considering any sampling approaches. In many studies [14], [15], [16], the classification of network attacks from an imbalanced binary class distribution has been looked at without taking sampling methods into account. Binary classification refers to classification problems where there are only two target classes. Imbalanced binary classification problems occur when one class has many more training examples than the other. Typically, the normal traffic class has a majority, while the abnormal (under attack) traffic class has a significantly smaller minority. The research by [14] aimed to build a classifier to determine whether a Distributed Denial of Service (DDoS) attack occurs on the network. The study employs a range of classifiers, including Extreme Gradient Boosting (XGB), Support Vector Machine (SVM), Logistic Regression (LR), K-Nearest Neighbor (KNN), and Decision Tree (DT). Evaluation metrics such as F1-score, Precision, Recall, and Accuracy indicate XGBoost's strength as the top-performing classifier, achieving an accuracy of 98.24%. In another study

for DDoS attack detection, the authors of [15] applied Logistic Regression, K-Nearest Neighbour, Multi-layer Perceptron, and Decision Tree to investigate the best detection model. Notably, KNN and DT demonstrate superior accuracy, especially for TCP and ICMP flooding attacks, while for UDP, DT exhibits a better accuracy of 77.23% with an almost equivalent F1-score.

Concurrently, there exists a group of researchers actively addressing the challenges associated with multiclass imbalanced scenarios in network attack classification. Examples of instances include [17], where the F1-score remains suboptimal, indicating the model is not achieving adequate performance on the minority class, even though overall accuracy appears high. In a study employing the UNSW-NB15 dataset [11], even though the dataset is multiclass imbalanced, the primary emphasis lies on presenting overall performance rather than individual class results. The findings demonstrate that Random Forest attains the best Area Under the Curve (AUC) and F2 scores. Additionally, [18] utilizes the NSL-KDD dataset, comparing the performance of Naive Bayes and SVM. Despite SVM's accuracy exceeding 90%, the F1-score remains around 0.69.

An additional study in [17] developed a model to classify benign network traffic versus malicious attack categories like Distributed Denial of Service (DDoS) attacks that leverage malicious TCP ACK or PSH-ACK packet flows.The results highlight the superiority of logistic regression over other classifiers used in the paper. The study in [19] applied the CICIDS2017 network intrusion detection benchmark to assess an array of both classical (Decision Tree, K-nearest Neighbours, and Support Vector Machine) and ensemble classifiers (Random Forest, GraBoost, and AdaBoost) for identifying malicious network behaviours within realistic traffic. The study reported that GraBoost outperformed other classifiers in terms of accuracy, precision, recall, and F1-score. Meanwhile, AdaBoost struggles with dataset complexity, lagging other classifiers significantly across all metrics.

Network security operates in a dynamic realm where cybersecurity threats continually evolve in complexity and diversity. The deployed classifier must constantly adapt to new attacks. However, a notable proportion of cybersecurity research concentrates on the development of machine learning models without considering the accurate detection of new attacks with a small sample size (minority classes). While studies like [14] and [15] offer insights into algorithmic performance in binary contexts, there exists a significant gap in understanding whether these algorithms retain their effectiveness amid the complexities of multiclass imbalanced datasets. Moreover, the interaction between different algorithms and metrics, such as the F1-score, remains underexplored. Therefore, a comprehensive investigation is needed to identify the high-performance algorithm that overcomes the imbalanced class distribution in the absence of sampling methods to rebalance the distribution. Additionally, the limited reporting of individual class results, as observed in [11], poses a gap in our understanding of algorithmic vulnerabilities and strengths across diverse attack types. Lastly, despite extensive algorithm testing, a systematic exploration of the suitability of different machine learning algorithm families for multiclass imbalanced datasets is lacking. Addressing these research gaps is imperative for advancing

the field, guiding algorithm selection, and advancing network intrusion detection in complex, real-world scenarios.

## III. METHODOLOGY

A series of experiments followed the procedure outlined in Fig. 2. Initially, the dataset was partitioned into two segments for training and testing purposes. Subsequent experiments evaluated the performance of four distinct traditional machine learning algorithms to identify the optimal base algorithm for the ensemble. Upon determining the optimal conventional algorithm, further tests were conducted to ascertain the most effective ensemble method, utilizing the previously selected traditional algorithm.



Fig. 2. Experimental evaluation flow.

### A. Dataset

In this research, we strategically utilize a highly imbalanced network intrusion dataset, reflective of real-world network anomaly scenarios, as our primary training resource. The dataset selected for this study is the publicly accessible and extensively recognized UNSW-NB15 dataset. It comprises ten different attack categories, represented by 43 features as detailed in Table I. This dataset includes a total of 257,673 instances, categorized into ten distinct classes, as delineated in Table II.

Table II highlights a notable characteristic of the UNSW-NB15 dataset: its classification as a multiclass imbalanced dataset. There are substantial variations in the frequency of different attack categories. These differences mirror the complexity of real-world scenarios, where certain network attacks, although less frequent, may be of higher significance. Categories such as Analysis, Backdoor, Reconnaissance, Shellcode, and Worms, each accounting for less than 6% of the total instances, are thus identified as minority classes in this study.

In order to demonstrate the skewed and highly imbalanced class scenario that exists within this dataset, we present the

TABLE I. FEATURES OF THE UNSW-NB15 DATASET

| No. | Features | Data types | No. | Features | Data types |
|---|---|---|---|---|---|
| 1 | id | int64 | 23 | dtrcpb | int64 |
| 2 | dur | float64 | 24 | dwin | int64 |
| 3 | proto | object | 25 | tcprtt | float64 |
| 4 | service | object | 26 | synack | float64 |
| 5 | state | object | 27 | ackdat | float64 |
| 6 | spkts | int64 | 28 | smean | int64 |
| 7 | dpkts | int64 | 29 | dmean | int64 |
| 8 | sbytes | int64 | 30 | trans-depth | int64 |
| 9 | dbytes | int64 | 31 | response-body-len | int64 |
| 10 | rate | float64 | 32 | ct-srv-src | int64 |
| 11 | sttl | int64 | 33 | ct-state-ttl | int64 |
| 12 | dttl | int64 | 34 | ct-dst-ltm | int64 |
| 13 | sload | float64 | 35 | ct-src-dport-ltm | int64 |
| 14 | dload | float64 | 36 | ct-dst-sport-ltm | int64 |
| 15 | sloss | int64 | 37 | ct-dst-src-ltm | int64 |
| 16 | dloss | int64 | 38 | is-ftp-login | int64 |
| 17 | sinpkt | float64 | 39 | ct-ftp-cmd | int64 |
| 18 | dinkpt | float64 | 40 | ct-flw-http-mthd | int64 |
| 19 | sjit | float64 | 41 | ct-src-ltm | int64 |
| 20 | djit | float64 | 42 | ct-src-dst | int64 |
| 21 | swin | int64 | 43 | is-sm-ips-ports | int64 |
| 22 | stcpb | int64 | 44 | attack-cat | object |

disparity between classes by utilising two metrics that are distinct from one another but interconnected metrics. Firstly, the Fraction to Majority Class was calculated using Eq. (1) as shown below:

$$\text{Fraction to Majority Class (\%)} = \frac{TNIPC}{TNIMMC} \times 100 \quad (1)$$

This metric aligns with the challenges identified in the practical scenario of network intrusion detection. Under these circumstances, some classes may have a low occurrence rate yet present a substantial risk. Eq. (2) was applied to calculate the Fraction to Total Instances is shown below:

$$\text{Fraction to Total Instances (\%)} = \frac{TNIPC}{TNIWD} \times 100 \quad (2)$$

For both Eq. (1) and Eq. (2), $TNIPC$ represents the total number of instances for a specific class, $TNIMMC$ is the total number of instances for the most majority class (the class with the highest number of instances), and $TNIWD$ indicates the total number of instances for the whole dataset. Instead of being mere mathematical equations, these equations also provide a clear understanding of the complex, imbalanced distribution of the dataset, which is also the problem found in the real-world situation.

By closely analyzing the imbalance and complexity of the dataset, a strong understanding of the complications of the dataset is established. This is crucial as it will ensure the techniques to be used are able to be utilized accurately and correctly when facing the imbalanced problem.

### B. Data Preparation

Before initiating the model training process, several preparatory steps are essential for the UNSW-NB15 dataset to ready the classifiers for subsequent stages. As indicated in Table I, the datatypes of the attack classes were initially in an object format. Consequently, the initial step in this research was to assign a numerical value to each class. This transformation is crucial as it not only standardizes representations but

TABLE II. NUMBER OF INSTANCES IN EACH ATTACK CLASS IN THE
UNSW-NB15 DATASET

| Classes (attack-cat) | Total number of instances | Fraction to Majority class (Percentage,%) | Fraction to Total instances (Percentage,%) |
|---|---|---|---|
| Analysis | 2,677 | 2.9 | 1.0 |
| Backdoor | 2,329 | 2.5 | 0.9 |
| DoS | 16,353 | 17.6 | 6.3 |
| Exploits | 44,525 | 48.9 | 17.3 |
| Fuzzers | 24,246 | 26.1 | 9.4 |
| Generic | 58,871 | 63.3 | 22.8 |
| Normal | 93,000 | 100 | 36.1 |
| Reconnaissance | 13,987 | 15.0 | 5.4 |
| Shellcode | 1,511 | 1.6 | 0.6 |
| Worms | 174 | 0.2 | 0.1 |
| **Total** | **257673** | | |

TABLE III. DATASET DISTRIBUTION

| Classes (attack-cat) | Assigned Number | Total number of Instances in UNSW-NB15 dataset | Number of Instances in Training dataset | Number of Instances in Testing dataset |
|---|---|---|---|---|
| Analysis | 0 | 2,677 | 1,874 | 803 |
| Backdoor | 1 | 2,329 | 1,630 | 699 |
| DoS | 2 | 16,353 | 11,447 | 4,906 |
| Exploits | 3 | 44,525 | 31,167 | 13,358 |
| Fuzzers | 4 | 24,246 | 16,972 | 7,274 |
| Generic | 5 | 58,871 | 41,210 | 17,661 |
| Normal | 6 | 93,000 | 65,100 | 27,900 |
| Reconnaissance | 7 | 13,987 | 9,791 | 4,196 |
| Shellcode | 8 | 1511 | 1,058 | 453 |
| Worms | 9 | 174 | 122 | 52 |
| **Total** | | 257,673 | 180,371 | 77,302 |

also ensures compatibility with machine learning algorithms. Furthermore, features such as proto, service, and state, which are initially in an object format, have also been encoded.

After assigning numerical values to the classes, the dataset underwent a stratified 70:30 split. Seventy percent of the data was allocated as the training dataset, enabling the classifier to learn patterns and relationships within the data. The remaining 30% served as the testing dataset, used to evaluate the performance of the trained classifiers in this research. The stratified split ensures equitable representation of all classes in both training and testing datasets, preventing any class from being overrepresented and potentially misleading classifier performance.

The detailed composition of the dataset split is presented in Table III. Employing the aforementioned stratified 70:30 split, instances for each class were proportionately divided between the training and testing datasets. This approach provides a more equitable and accurate assessment of the performance of the algorithms used in this research, particularly in addressing multiclass imbalanced classification challenges.

### C. Conventional Machine Learning Algorithms

This paper evaluates four distinct conventional machine learning algorithms, each representing a different family of algorithms: tree-based, instance-based, function-based, and Bayesian-based. These algorithms were chosen for their simplicity and computational efficiency, a desirable trait given the need for rapid training in scenarios involving frequently

updated network attacks. The assessed algorithms are: Decision Tree (DT) from the tree-based family, K-nearest neighbor (KNN) from the instance-based family, Multilayer Perceptron (MLP) from the function-based family, and Naive Bayes (NB) from the Bayesian-based family. Initially, the performance of these algorithms is evaluated to identify the most effective family-based classifier for addressing the multiclass imbalanced problem. A brief description of these algorithms is as follows:

1) Decision Tree (DT): A well-known approach used in the field of network intrusion detection. It constructs a hierarchical tree with decision leaves and data element nodes to solve the classification problem. Although [20] has raised concerns about the necessity for numerous splits in a skewed distribution dataset, some researchers [21], [22], [23] have proved the efficiency of DT in this field.

2) K-nearest neighbor (KNN): An instance-based algorithm, KNN classifies dataset instances using Euclidean distance to measure the proximity between training and testing instances [24]. It is simple and robust against noisy data [25], albeit with some efficiency drawbacks, particularly in selecting the optimal "$k$" value [26]. In our experiment, $k = 10$ was chosen as the most suitable value after fine-tuning.

3) Multilayer Perceptron (MLP): As a neural network, or function-based algorithm, MLP consists of multiple interconnected neuron layers [27], [26]. The number of hidden and output layers determines its structure [28]. In our experiments, we configured the MLP with 100 hidden layers, using the Rectified Linear Unit (ReLU) as the activation function and Adam as the optimizer with a learning rate of 0.001. The maximum number of iterations was set to 200.

4) Naive Bayes (NB): Naive Bayes classifier is a family of simple probabilistic classification algorithms based on Bayes' theorem. In contrast to Bayes theorem, it is designed based on naive assumption that features are independent from each other to simplify the algorithm. In this experiment, we implemented Gaussian variant which uses Gaussian Distribution for the feature values of each class [29]. Instead of solely relying on the Euclidean distance from the class mean, this algorithm takes both into account. Yet, it does have the drawback of only modeling each dimension independently, as it neglects the joint distribution of weight and height [30].

### D. Tree-based Ensemble Algorithms

The research employed a selected set of ensemble algorithms, with a specific focus on tree-based families in which the decision tree serves as the primary classifier for these methods. The choice was made due to the decision tree's ability to handle the imbalanced dataset, as was discussed in Section V-A.

The following nine tree-based ensemble algorithms were applied in this study:

1) Bagging: Bagging (Bootstrap Aggregating) is an effective technique that is able to solve the high

variance problem faced by some algorithms, such as decision trees. It involves constructing several trees without pruning and is able to show reliable results [31].

2) Random Forest (RF): An improved version of Bagging that is able to reduce noise, solve outliers, and overfit problems, which are common challenges found in a dataset. By reducing correlations between individual classifiers, RF effectively eliminates and deals with these difficulties, creating a robust and reliable model [31], [32].

3) Extremely Randomized Tree (ExtraTree): As compared to RF, this algorithm, which is also an evolution of Bagging, constructs random trees by using the instances of the dataset [31]. An enhanced robustness and increased resilience were able to be guaranteed with this intentional injection of diversity, which also strengthened the overall ensemble.

4) Adaptive Boosting (AdaBoost): It is a weighted-assigned ensemble algorithm that modifies the weight of the instances of the dataset dynamically. By doing so, the algorithm is able to allocate attention strategically during the construction of subsequent models, which enhances its capabilities in handling the different complexities present in the network intrusion data.

5) Gradient Boosting (GraBoost): GraBoost is a very complex and sophisticated ensemble algorithm. Despite its complexity, GraBoost stands as one of the most formidable ensemble methods, particularly distinguished for its efficacy in elevating classification performance amidst the challenges posed by imbalanced datasets [31].

6) Histogram-based Gradient Boosting (HistGraBoost): HistGraBoost, an innovative boosting algorithm, addresses a key limitation of the GB algorithm—lengthy training times on large datasets. This is remedied by discretizing continuous input variables, optimizing efficiency. The critical hyperparameter is the learning rate, with extensive optimization through multiple rounds of tuning [33].

7) Extreme Gradient Boosting (XGBoost): XGBoost, a highly scalable tree boosting system, is renowned for state-of-the-art performance in machine learning challenges. Leveraging sparsity-aware techniques and insights into cache access patterns, data compression, and sharding, XGBoost excels in efficiency. It outperforms comparable systems on large datasets while optimizing resource utilization [34].

8) Light Gradient Boosting Machine (LightGBM): LightGBM, a robust framework implementing Gradient Boosted Decison Tree (GBDT), emphasizes efficient parallel training. With features like accelerated training speed, reduced memory consumption, and support for distribution, LightGBM excels in accuracy and swift processing of massive datasets [35].

9) Category Gradient Boosting (CatBoost): CatBoost, an innovative algorithm, automatically treats categorical features as numerical characteristics. Utilizing a combination of category features enriches feature dimensions, while a perfectly symmetrical tree model reduces overfitting, enhancing accuracy and generalizability [36]. This categorical-centric approach positions CatBoost as a sophisticated solution for handling categorical features within gradient boosting algorithms.

The strategic evaluation of these ensemble algorithms is necessary to tackle the intricate challenges posed by imbalanced datasets. Section V-A will showcase the preliminary outcomes that demonstrate the proficiency of each traditional machine learning algorithm. This will provide an understanding of the factor influences the selection of algorithms in this research project.

## IV. EVALUATION METRICS

The F1-score is crucial in evaluating the effectiveness of the tree-based ensemble methods used in this work. The F1-score is instrumental in situations where imbalances are common. It offers a balanced evaluation that considers the constraints of accuracy, which can often give too much weight to classes with a high number of instances or overlook differences within classes. The decision to prioritize the F1-score as the primary evaluation metric is based on its inherent insensitivity to class imbalance. It is a suitable tool for assuring an unbiased and impartial assessment [37].

The F1-score is mathematically defined in Eq.(3).

$$\text{F1-score} = \frac{2 \cdot precision \cdot recall}{precision + recall} \qquad (3)$$

where precision refers to the measure of how accurate positive predictions are. It is calculated by dividing the number of true positive ($TP$) by the sum of true positive and false positive ($FP$) predictions. Recall, also known as sensitivity or the true positive rate, gauges the ability to accurately identify positive instances, measured as the ratio of true positive ($TP$) predictions to the sum of true positive and false negative ($FN$) predictions.

The F1-score ranges between $0$ and $1$, with $1$ denoting optimal performance. A higher F1-score signifies superior performance, achieving an equilibrium between precision and recall [38]. This paper also reports the F1-score for each class to provide insights into class-specific performance. Understanding how well the model performs for each class is essential in real-world applications. Beyond complementing the F1-score per class, we also provided the Weighted F1-score and the Macro Average F1-score to analyze the overall performance of the algorithms.

The macro average F1-score assigns equal importance to each class, so preventing the dominance of larger classes from overshadowing the performance of smaller ones. Additionally, it offers valuable insights into the performance of each class separately, which is particularly useful in situations when the performance of each class is of utmost importance. The weighted F1-score enables the allocation of distinct weights to classes according to their significance, so effectively addressing imbalances in a manner better to macro averaging. In this experiment, the weights are allocated according to the sizes of the classes. Note that this research excludes the use of micro average F1-score due to its susceptibility to being

influenced by classes with bigger sizes, which may result in the performance of smaller classes being disregarded.

The equation for the Weighted F1-score is provided in Eq. (4), whereas the equation for the Macro Average F1-score is given in Eq. (5).

$$\text{Weighted F1-score} = \frac{\sum_i \text{F1-score}_i \times \text{Weight}_i}{\sum_i \text{Weight}_i} \qquad (4)$$

where F1-score$_i$ is the F1-score for class $i$, and Weight$_i$ is the weight assigned to class $i$ which refers to the proportion of instances in class $i$ in the dataset.

$$\text{Macro Average F1-score} = \frac{\sum_{i=1}^{N} \text{F1-score}_i}{N} \qquad (5)$$

where F1-score$_i$ is the F1-score for class $i$ and $N$ is the number of classes.

In addition to the F1-score, we also deliver the results based on the accuracy value.

$$Accuracy(\%) = \frac{TP + TN}{TP + TN + FP + FN} \times 100 \qquad (6)$$

where $TP + TN$ denotes the total number of instances correctly classified in that class, and $TP + TN + FP + FN$ represent the total number of instances in that class in the testing dataset.

In summary, the chosen evaluation measures, with the F1-score as the leading indicator, provide a thorough and informative insight for evaluating the effectiveness of the measured algorithms on imbalanced multiclass datasets. The F1-score enables us to assess the efficacy better. The method aimed to improve classification performance, especially for rare classes in real-world situations.

## V. RESULTS AND DISCUSSION

### A. Preliminary Results

Four machine learning methods from different families were utilized to train the classifier on the training dataset and then test it on the testing dataset. The algorithmic selection consisted of representatives from various families, including Decision Tree (DT) from the tree-based family, K-nearest neighbour (KNN) from the instance-based family, Multilayer Perceptron (MLP) from the functions-based family, and Naive Bayes (NB) from the Bayesian-based family. The purpose of this selection process was to identify the most suitable machine learning algorithm families for tackling the complex task of multiclass imbalanced classification.

The thorough assessment, as depicted in Tables IV and V, showcases the results of our study. The performance of the machine learning algorithms varies considerably across different attack categories. The Decision Tree (DT) approach demonstrated the maximum accuracy in the "Generic" class with 98.31% and the "Normal" class with 91.26%. Nevertheless, the Multilayer Perceptron (MLP) exhibited higher accuracy in the "Normal" class with 99.60%. When working with

classes that have a small number of instances, like "Worms" and "Shellcode," even a single misclassification might have a large impact on the accuracy results due to the low size of the sample. The results show that in overall, the MLP exhibits inferior accuracy compared to other algorithms, indicating that it may encounter difficulties handling the complex nature of specific attack patterns. The mean accuracy for DT is 48.55%; for KNN, it is 31.57%; for MLP, it is 3.08%; and for NB, it is 16.76% across all attack classes. These results provide a perspective on the performance of algorithms. Still, the interpretation should be done carefully, considering the presence of class imbalances.

The tree-based classifiers, specifically the Decision Tree (DT) with the highest Weighted F1-score of 0.80, clearly outperformed the algorithms from other families regarding overall F1-scores and accuracy. The Weighted F1-score of KNN is 0.65, which is the second highest among the models. Naive Bayes is entirely ineffective in detecting Analysis and Denial of Service (DoS) threats. The KNN, MLP and NB were facing difficulties in accurately detecting and categorizing threats such as Analysis, Backdoors, Shellcode, and Worms as the F1-score for each class is below 0.15. The MLP exhibited poor results with all classes except for the "Normal" class, achieving an F1-score of 0.1 or lower. This demonstrates that the MLP is only capable of recognizing regular network traffic and lacks the ability to identify network attacks.

The experiment strongly suggests a greater efficacy of the decision tree, as evidenced by the substantial findings. Furthermore, a per-class analysis reveals that it surpassed other traditional algorithms in performance for all classes. This discovery yields a vital inference: tree-based algorithms demonstrate superior performance when addressing multiclass imbalanced classification issues compared to conventional techniques. This analysis clarifies the reasoning for choosing tree-based ensemble techniques and explores the further findings.

### B. The Evaluation of Tree-based Ensemble Algorithms Performance

In this part, we will further investigate the most appropriate tree-based technique for practical application in the problem of multiclass imbalanced classification. This analysis is based on the findings presented in Section V-A and focuses on comparing different tree-based ensemble algorithms. This section offers an extended analysis of the tree-based ensemble algorithms employed in this study. Similar to the previous section (Section V-A), the selected tree-based ensemble algorithms were evaluated based on the accuracy, F1-score per class, Weighted F1-score and Macro Average F1-score as explained in Section IV.

The tables labelled as VI and VII include valuable information about how well these ensemble approaches, built on trees, perform in classifying instances for each category. XGBoost outperforms other algorithms, demonstrating exceptional results across the majority of classes with a classification accuracy of 50%. However, it is essential to note that there are outliers within the Analysis, Backdoor, and Denial of Service (DoS) attack categories. All the algorithms used in this study exhibit reduced accuracy in those instances.

TABLE IV. ACCURACY RESULTS FOR FOUR CONVENTIONAL MACHINE LEARNING ALGORITHMS FOR EACH ATTACK CLASS IN UNSW-NB15 DATASET

| Attack class | Number of instances per classes in test dataset | Accuracy (%) | | | |
|---|---|---|---|---|---|
| | | Decision Tree(DT) | K-nearest neighbor(KNN) | Multilayer Perceptron(MLP) | Naive Bayes(NB) |
| Analysis | 803 | **12.07** | 2.74 | 0.62 | 0.00 |
| Backdoor | 699 | **9.16** | 0.14 | 0.29 | 0.29 |
| DoS | 4,906 | **33.87** | 30.54 | 6.38 | 0.18 |
| Exploits | 13,358 | **73.87** | 50.94 | 0.91 | 2.00 |
| Fuzzers | 7,274 | **58.62** | 26.23 | 1.94 | 21.37 |
| Generic | 17,661 | **98.31** | 97.42 | 0.05 | 97.99 |
| Normal | 27,900 | 91.26 | 79.73 | **99.60** | 45.67 |
| Reconnaissance | 4,196 | **75.92** | 33.92 | 0.00 | 43.07 |
| Shellcode | 453 | **59.02** | 7.95 | 0.00 | 1.33 |
| Worms | 52 | **53.85** | 0.00 | 0.00 | 3.85 |
| **Average** | | **48.55** | 31.57 | 3.08 | 16.76 |

TABLE V. F1-SCORE RESULTS FOR FOUR CONVENTIONAL MACHINE LEARNING ALGORITHMS FOR EACH ATTACK CLASS IN UNSW-NB15 DATASET

| Attack class | Number of instances per classes in test dataset | F1-scores | | | |
|---|---|---|---|---|---|
| | | Decision Tree(DT) | K-nearest neighbor(KNN) | Multilayer Perceptron(MLP) | Naive Bayes(NB) |
| Analysis | 803 | **0.17** | 0.05 | 0.01 | 0.00 |
| Backdoor | 699 | **0.15** | 0.00 | 0.01 | 0.01 |
| DoS | 4,906 | **0.33** | 0.30 | 0.10 | 0.00 |
| Exploits | 13,358 | **0.69** | 0.43 | 0.02 | 0.04 |
| Fuzzers | 7,274 | **0.61** | 0.30 | 0.04 | 0.25 |
| Generic | 17,661 | **0.99** | 0.98 | 0 | 0.64 |
| Normal | 27,900 | **0.91** | 0.78 | 0.54 | 0.61 |
| Reconnaissance | 4,196 | **0.82** | 0.49 | 0.00 | 0.15 |
| Shellcode | 453 | **0.59** | 0.12 | 0.00 | 0.02 |
| Worms | 52 | **0.47** | 0.00 | 0.00 | 0.01 |
| **Macro Average F1-score** | | **0.63** | 0.31 | 0.06 | 0.16 |
| **Weighted F1-score** | | **0.80** | 0.65 | 0.21 | 0.40 |

For the Analysis class, RF, XGBoost, and Bagging perform better than other models, attaining the highest F1-scores of 0.19. The F1-scores for DT and ExtraTree are both 0.17. XGBoost outperforms other models in terms of F1-score with a value of 0.17 for the Backdoor class. Bagging, Decision Trees (DT), and Random Forest (RF) exhibit similar performance, with F1-scores almost equal to 0.16. ExtraTree and DT demonstrate the most robust performance in terms of F1-score (0.33) for the DoS attack. Bagging and Random Forest (RF) perform strongly, achieving F1-scores ranging from 0.30 to 0.33. We found that XGBoost and CatBoost achieved the highest F1-score of 0.74 in the Exploit class. Other methods such as Bagging, RF, ExtraTree, GraBoost, and HistGraBoost produce comparable scores ranging from 0.72 to 0.73. For Fuzzer attacks, Bagging demonstrates the most significant F1-scores, precisely 0.66. Other algorithms, such as Random Forest (RF) and ExtraTree, attain scores about equal to 0.65. The results also show that most algorithms perform exceptionally in categorizing generic traffic, as evidenced by their high F1-scores of approximately 0.99. For Normal traffic, the results show that XGBoost, Bagging, RF, and ExtraTrees exhibit the most outstanding F1-scores of 0.93. Bagging and XGBoost achieve the highest F1-scores for the Reconnaissance and Shellcode classes, with 0.84 and 0.69, respectively. Regarding the class with the smallest number of samples, Worms, it is observed that XGBoost attains the highest F1-scores, precisely 0.63.

The F1-score data shown in Table VI demonstrates that advanced ensemble approaches, namely Bagging, Random Forest, XGBoost, and ExtraTrees, exhibited superior performance compared to the conventional Decision Tree. Bagging, Random Forest, XGBoost, and ExtraTree obtained the highest Weighted F1-score. Bagging attains its highest macro average F1-score of 0.63, denoting outstanding overall performance. Additional algorithms, such as XGBoost and DT, exhibit similar performance with macro F1-scores ranging from 0.60 to 0.63. Bagging, Random Forest, ExtraTrees, and XGBoost demonstrate superior performance in addressing imbalanced classes, as evidenced by their highest weighted F1-scores of 0.82. Several other algorithms, such as DT, GraBoost, HistGraBoost, LightGBM, and CatBoost, achieve weighted F1-scores in the range of 0.79−0.80. The weighted F1-score offers a more accurate evaluation by taking into account both the performance of each class and the distribution of classes.

The results suggest that the algorithm's efficacy is significantly influenced by the distinctive properties of each class, thereby necessitating an understanding of attack characteristics. After analyzing Table VII, it is evident that XGBoost emerges as the most robust choice, outperforming all other tree-based ensemble algorithms by attaining the highest F1-score for eight out of ten classes. While Bagging demonstrates comparable Weighted F1-scores and Macro Average F1-scores, an in-depth analysis indicates that XGBoost surpasses in particular categories (except for DoS and Fuzzers). Furthermore, the accuracy results presented in Table VIII conclusively indicate that XGBoost exceeded other algorithms in terms of accuracy, with Bagging being an equally strong candidate.

The data presented in this paper suggest that XGBoost and Bagging is the best tree-based ensemble method for multiclass imbalanced classification in the particular scenario of network attack detection. The study's results emphasize the algorithm's effectiveness in tackling the difficulties posed by imbalanced datasets, making it a highly appropriate choice for practical

TABLE VI. F1-SCORE RESULTS FOR DECISION TREE AND NINE TREE-BASED ENSEMBLE MACHINE LEARNING ALGORITHMS FOR EACH ATTACK CLASS IN THE UNSW-NB15 DATASET

| Attack class | Number of instances per classes in test dataset | F1-scores | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | DT | Bagging | RF | ExtraTree | AdaBoost | GraBoost | HistGraBoost | XGBoost | LightGBM | CatBoost |
| Analysis | 803 | 0.17 | **0.19** | **0.19** | 0.17 | 0.11 | 0.08 | 0.13 | **0.19** | 0.13 | 0.09 |
| Backdoor | 699 | 0.16 | 0.16 | 0.16 | 0.15 | 0.01 | 0.12 | 0.14 | **0.17** | 0.09 | 0.12 |
| DoS | 4,906 | **0.33** | 0.32 | 0.30 | **0.33** | 0.08 | 0.12 | 0.08 | 0.20 | 0.25 | 0.17 |
| Exploits | 13,358 | 0.69 | 0.72 | 0.72 | 0.72 | 0.22 | 0.72 | 0.73 | **0.74** | 0.70 | **0.74** |
| Fuzzers | 7,274 | 0.61 | **0.66** | 0.65 | 0.65 | 0.29 | 0.58 | 0.52 | 0.64 | 0.59 | 0.61 |
| Generic | 17,661 | **0.99** | **0.99** | **0.99** | **0.99** | 0.93 | **0.99** | **0.99** | **0.99** | **0.99** | **0.99** |
| Normal | 27,900 | 0.91 | **0.93** | **0.93** | **0.93** | 0.60 | 0.91 | 0.91 | **0.93** | 0.90 | 0.92 |
| Reconnaissance | 4,196 | 0.82 | **0.84** | 0.83 | 0.82 | 0.51 | 0.83 | 0.83 | **0.84** | 0.80 | 0.83 |
| Shellcode | 453 | 0.61 | **0.69** | 0.67 | 0.64 | 0.39 | 0.65 | 0.53 | **0.69** | 0.52 | 0.61 |
| Worms | 52 | 0.47 | 0.62 | 0.20 | 0.20 | 0.01 | 0.17 | 0.14 | **0.63** | 0.13 | 0.17 |
| **Macro Average F1-score** | | 0.60 | **0.63** | 0.59 | 0.58 | 0.33 | 0.54 | 0.53 | 0.62 | 0.54 | 0.55 |
| **Weighted F1-score** | | 0.80 | **0.82** | 0.82 | 0.82 | 0.53 | 0.79 | 0.79 | **0.82** | 0.79 | 0.80 |

TABLE VII. ALGORITHMS WITH HIGHEST F1-SCORE PER CLASS

| Class | DT | Bagging | RF | ExtraTree | AdaBoost | GraBoost | HistGraBoost | XGBoost | LightGBM | CatBoost |
|---|---|---|---|---|---|---|---|---|---|---|
| Analysis | | ✓(0.19) | ✓(0.19) | | | | | ✓(0.19) | | |
| Backdoor | | | | | | | | ✓(0.17) | | |
| DoS | ✓(0.33) | | | ✓(0.33) | | | | | | |
| Exploits | | | | | | | | ✓(0.74) | | ✓(0.74) |
| Fuzzers | | ✓(0.66) | | | | | | | | |
| Generic | ✓(0.99) | ✓(0.99) | ✓(0.99) | ✓(0.99) | | ✓(0.99) | ✓(0.99) | ✓(0.99) | | ✓(0.99) |
| Normal | | ✓(0.93) | ✓(0.93) | ✓(0.93) | | | | ✓(0.93) | | |
| Reconnaissance | | ✓(0.84) | | | | | | ✓(0.84) | | ✓(0.83) |
| Shellcode | | ✓(0.69) | | | | | | ✓(0.69) | | |
| Worms | | | | | | | | ✓(0.63) | | |

implementation in cybersecurity and network intrusion detection.

## VI. CONCLUSION AND FUTURE WORKS

The findings indicate that tree-based ensemble methods, including Bagging, Random Forest, XGBoost, and ExtraTrees, have achieved a high Weighted F1-score, despite the constraint of an imbalanced training dataset. These qualities make them very suitable for identifying network intrusions in the UNSW-NB15 dataset. XGBoost surpassses other tree-based algorithms in terms of per-class F1-scores, which is a useful performance measure for addressing multiclass imbalance problems. Nevertheless, the overall accuracy of XGBoost is about equivalent to that of Bagging. These findings confirm that XGBoost is the most effective approach for addressing multiclass imbalance classification, with Bagging being the most viable option. In summary, the results highlight the effectiveness of Decision Tree (DT) and tree-based ensemble algorithms in handling the problem of imbalanced multi-class datasets.

This study has offered valuable insights into the efficacy of tree-based ensemble algorithms for multiclass imbalanced classification in network intrusion detection. However, it is crucial to recognise the underlying constraints and difficulties. Although ensemble strategies have been used to address class imbalance, the problem persists. The disproportionate allocation of classes, specifically pertaining to minority categories such as Analysis, Backdoor, and Denial of Service (DoS), still poses substantial difficulties in detecting these classes.

Beyond the difficulties and constraints, the outcomes provide a solid groundwork for future studies in this domain. Further investigations into feature engineering, advanced sampling approaches, or algorithmic adaptations that can effectively improve the identification of minority class occurrences

should be conducted. More advanced algorithms with adaptive sampling capable of dealing with data changes over time are likely needed. These efforts are necessary for creating resilient and flexible solutions to address the constantly evolving cyber-attack scenario.

Future research must consider developing and using domain-specific evaluation metrics that improve the interpretation of algorithmic performance in situations with imbalances across many classes. This evaluation metric should surpass conventional performance metrics such as the F1-score. It should provide a comprehensive assessment considering the trade-off between false positives and false negatives in various domains. This will result in a more thorough evaluation of performance.

## ACKNOWLEDGMENT

## REFERENCES

[1] I. Lella, E. Tsekmezoglou, M. Theocharidou, E. Magonara, A. Malatras, R. S. Naydenov, and C. Ciobanu, "Enisa threat landscape 2023," 2023.

[2] M. Ring, S. Wunderlich, D. Scheuring, D. Landes, and A. Hotho, "A survey of network-based intrusion detection data sets," *Computers & Security*, vol. 86, pp. 147–167, 2019.

[3] K. M. Hasib, M. S. Iqbal, F. M. Shah, J. Al Mahmud, M. H. Popel, M. I. H. Showrov, S. Ahmed, and O. Rahman, "A survey of methods for managing the classification and solution of data imbalance problem," *Journal of Computer Science*, vol. 16, p. 1546–1557, Nov. 2020.

[4] S. K. M, Devidas, S. N. Pai, S. Kolekar, V. Pai, and B. R, "Use of machine learning and random oversampling in stroke prediction," in *2022 International Conference on Artificial Intelligence and Data Engineering (AIDE)*, pp. 331–337, 2022.

TABLE VIII. ACCURACY RESULTS FOR DECISION TREE AND NINE TREE-BASED ENSEMBLE MACHINE LEARNING ALGORITHMS FOR EACH ATTACK CLASS IN THE UNSW-NB15 DATASET

| Attack class | Number of instances per classes in test dataset | Accuracy (%) | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | DT | Bagging | RF | ExtraTree | AdaBoost | GraBoost | HistGraBoost | XGBoost | LightGBM | CatBoost |
| Analysis | 803 | 12.08 | 11.08 | 10.96 | 10.84 | **62.15** | 4.36 | 7.11 | 10.35 | 8.48 | 4.99 |
| Backdoor | 699 | **9.73** | 9.30 | 9.01 | 8.87 | 4.44 | 6.29 | 8.01 | 9.16 | 5.58 | 6.72 |
| DoS | 4,906 | **33.93** | 28.52 | 26.17 | 31.60 | 4.84 | 6.84 | 4.33 | 12.87 | 18.23 | 10.43 |
| Exploits | 13,358 | 73.84 | 81.15 | 81.51 | 79.46 | 14.59 | 90.38 | **92.17** | 90.70 | 83.88 | 90.65 |
| Fuzzers | 7,274 | 58.48 | **61.48** | 60.69 | 60.16 | 21.10 | 50.67 | 40.17 | 58.09 | 52.95 | 57.85 |
| Generic | 17,661 | 98.23 | **98.33** | 97.97 | 97.94 | 89.78 | 97.67 | 97.93 | 98.28 | 97.98 | 97.94 |
| Normal | 27,900 | 91.32 | 94.87 | 94.44 | 94.55 | 51.80 | 92.66 | **95.75** | 95.10 | 90.20 | 94.14 |
| Reconnaissance | 4,196 | 76.03 | 76.34 | 76.23 | 75.13 | **79.62** | 76.77 | 75.89 | 76.94 | 76.44 | 76.11 |
| Shellcode | 453 | 61.16 | 70.87 | 66.46 | 60.27 | 29.15 | 65.79 | 54.98 | **73.08** | 56.08 | 59.83 |
| Worms | 52 | 57.70 | **67.32** | 13.46 | 13.46 | 59.63 | 48.09 | 38.47 | **67.32** | 32.70 | 9.63 |
| **Correctly Identified Instances** | | 62264 | 64251 | 63889 | 63793 | 38056 | 62812 | 63015 | **64598** | 61998 | 63959 |
| **Accuracy (%)** | | 80.55% | 83.12% | 82.65% | 82.52% | 49.23% | 81.25% | 81.52% | **83.57%** | 80.20% | 82.48% |

[5] F. Kamalov, H.-H. Leung, and A. K. Cherukuri, "Keep it simple: random oversampling for imbalanced data," in *2023 Advances in Science and Engineering Technology International Conferences (ASET)*, pp. 1–4, 2023.

[6] J. Hancock, T. M. Khoshgoftaar, and J. M. Johnson, "The effects of random undersampling for big data medicare fraud detection," in *2022 IEEE International Conference on Service-Oriented System Engineering (SOSE)*, pp. 141–146, 2022.

[7] Q. Dai, J.-w. Liu, and Y. Liu, "Multi-granularity relabeled under-sampling algorithm for imbalanced data," *Applied Soft Computing*, vol. 124, p. 109083, 2022.

[8] M. Muntasir Nishat, F. Faisal, I. Jahan Ratul, A. Al-Monsur, A. M. Ar-Rafi, S. M. Nasrullah, M. T. Reza, and M. R. H. Khan, "A comprehensive investigation of the performances of different machine learning classifiers with smote-enn oversampling technique and hyper-parameter optimization for imbalanced heart failure dataset," *Scientific Programming*, vol. 2022, pp. 1–17, 2022.

[9] S. Demir and E. K. Şahin, "Evaluation of oversampling methods (over, smote, and rose) in classifying soil liquefaction dataset based on svm, rf, and naïve bayes," *Avrupa Bilim ve Teknoloji Dergisi*, no. 34, pp. 142–147, 2022.

[10] S. Abokadr, A. Azman, H. Hamdan, and N. Amelina, "Handling imbalanced data for improved classification performance: Methods and challenges," in *2023 3rd International Conference on Emerging Smart Technologies and Applications (eSmarTA)*, pp. 1–8, 2023.

[11] I. Fosić, D. Žagar, and K. Grgić, "Network traffic verification based on a public dataset for ids systems and machine learning classification algorithms," in *2022 45th Jubilee International Convention on Information, Communication and Electronic Technology (MIPRO)*, pp. 1037–1041, 2022.

[12] N. Abedzadeh and M. Jacobs, "A survey in techniques for imbalanced intrusion detection system datasets," *International Journal of Computer and Systems Engineering*, vol. 17, no. 1, pp. 9 – 18, 2023.

[13] M. Kim and K.-B. Hwang, "An empirical evaluation of sampling methods for the classification of imbalanced data," *PLoS One*, vol. 17, no. 7, p. e0271260, 2022.

[14] R. Raj and S. Singh Kang, "Mitigating ddos attack using machine learning approach in sdn," in *2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICAC3N)*, pp. 462–467, 2022.

[15] P. S. Patil, S. L. Deshpande, G. S. Hukkeri, R. H. Goudar, and P. Siddarkar, "Prediction of ddos flooding attack using machine learning models," in *2022 Third International Conference on Smart Technologies in Computing, Electrical and Electronics (ICSTCEE)*, pp. 1–6, 2022.

[16] M. A. Talukder, K. F. Hasan, M. M. Islam, M. A. Uddin, A. Akhter, M. A. Yousuf, F. Alharbi, and M. A. Moni, "A dependable hybrid machine learning model for network intrusion detection," *Journal of Information Security and Applications*, vol. 72, p. 103405, 2023.

[17] S. K. Naing and T. T. Thwel, "A study of ddos attack classification using machine learning classifiers," in *2023 IEEE Conference on Computer Applications (ICCA)*, pp. 108–112, 2023.

[18] V. Santhi, J. Priyadharshini, M. Swetha, and K. Dhanavandhana, "A hybrid feature extraction method with machine learning for detecting the presence of network attacks," in *2023 International Conference on*

*Intelligent Systems for Communication, IoT and Security (ICISCoIS)*, pp. 454–459, 2023.

[19] R. Wen and K. Zhang, "Research on automated classification method of network attacking based on gradient boosting decision tree," in *2022 International Conference on Machine Learning and Knowledge Engineering (MLKE)*, pp. 72–76, 2022.

[20] F. Shakeel, A. S. Sabhitha, and S. Sharma, "Exploratory review on class imbalance problem: An overview," in *8th International Conference on Computer Communications and Networks Technologies (ICCCNT)*, 2017.

[21] N. Elmrabit, F. Zhou, F. Li, and H. Zhou, "Evaluation of machine learning algorithms for anomaly detection," in *International Conference on Cyber Security and Protection of Digital Services (Cyber Security)*, pp. 1–6, 2020.

[22] D. Kurniabudi, D. Stiawan, M. Y. Bin Bin Idris, A. M. Bamhdi, and R. Budiarto, "Improving the anomaly detection by combining pso search methods and j48 algorithm," *IEEE Explore for Emerging Cyber Security and Information Systems*, pp. 119–126, 2020.

[23] D. Kurniabudi, D. Stiawan, M. Y. Bin Bin Idris, A. M. Bamhdi, and R. Budiarto, "Cicids-2017 dataset feature analysis with information gain for anomaly detection," *IEEE Access*, vol. 8, pp. 132911–132921, 2020.

[24] N. Ali, D. Neagu, and P. Trundle, "Evaluation of k-nearest neighbour classifier performance for heterogeneous data sets," *SN Applied Sciences*, vol. 1, pp. 1–15, 2019.

[25] P. K. Syriopoulos, N. G. Kalampalikis, S. B. Kotsiantis, and M. N. Vrahatis, "k nn classification: a review," *Annals of Mathematics and Artificial Intelligence*, pp. 1–33, 2023.

[26] E. Y. Boateng, J. Otoo, and D. A. Abaye, "Basic tenets of classification algorithms k-nearest-neighbor, support vector machine, random forest and neural network: a review," *Journal of Data Analysis and Information Processing*, vol. 8, no. 4, pp. 341–357, 2020.

[27] H. Taud and J. Mas, "Multilayer perceptron (mlp)," *Geomatic approaches for modeling land change scenarios*, pp. 451–455, 2018.

[28] Q. Jiang, L. Zhu, C. Shu, and V. Sekar, "Multilayer perceptron neural network activated by adaptive gaussian radial basis function and its application to predict lid-driven cavity flow," *Acta Mechanica Sinica*, pp. 1–16, 2021.

[29] A. H. Jahromi and M. Taheri, "A non-parametric mixture of gaussian naive bayes classifiers based on local independent features," in *2017 Artificial Intelligence and Signal Processing Conference (AISP)*, pp. 209–212, 2017.

[30] R. D. Raizada and Y.-S. Lee, "Smoothness without smoothing: why gaussian naive bayes is not naive for multi-subject searchlight studies," *PLoS one*, vol. 8, no. 7, p. e69566, 2013.

[31] J. Brownlee, *Machine Learning Mastery with Python: Understand Your Data, Create Accurate Models and Work Projects End-To-End*. v1.19 ed., 2020.

[32] R. Saravanan and P. Sujatha, "A state of art techniques on machine learning algorithms: A perspective of supervised learning approaches in data classification," in *2019 2nd International Conference on Intelligent Computing and Control Systems (ICICCS)*, pp. 945–949, 2019.

[33] G. Marvin, L. Grbčić, S. Družeta, and L. Kranjčević, "Water distribution network leak localization with histogram-based gradient boosting," *Journal of Hydroinformatics*, vol. 25, no. 3, pp. 663–684, 2023.

[34] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, ACM, Aug. 2016.

[35] C. Ma, Y. Chi, D. Hao, and X. Ji, "A new approach based on feature selection of light gradient boosting machine and transformer to predict circrna-disease associations," *IEEE Access*, vol. 11, pp. 47187–47201, 2023.

[36] M. Luo, Y. Wang, Y. Xie, L. Zhou, J. Qiao, S. Qiu, and Y. Sun, "Combination of feature selection and catboost for prediction: The first application to the estimation of aboveground biomass," *Forests*, vol. 12, no. 2, p. 216, 2021.

[37] J. Wu, Z. Zhao, C. Sun, R. Yan, and X. Chen, "Learning from class-imbalanced data with a model-agnostic framework for machine intelligent diagnosis," *Reliability Engineering & System Safety*, vol. 216, p. 107934, December 2021.

[38] A. Géron, *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow.* " O'Reilly Media, Inc.", 2022.

# Bystander Detection: Automatic Labeling Techniques using Feature Selection and Machine Learning

Anamika Gupta[1], Khushboo Thakkar[2], Veenu Bhasin[3], Aman Tiwari[4], Vibhor Mathur[5]

S.S. College of Business Studies, University of Delhi, India[1,2,4,5]

P.G.D.A.V. College, University of Delhi, India[3]

*Abstract*—A hostile or aggressive behavior on an online platform by an individual or a group of people is termed as cyberbullying. A bystander is the one who sees or knows about such incidences of cyberbullying. A defender who intervenes can mitigate the impact of bullying, an instigator who accomplices the bully, can add to the victim's suffering, and an impartial onlooker who remains neutral and observes the scenario without getting engaged. Studying the behavior of Bystanders role can help in shaping the scale and progression of bullying incidents. However, the lack of data hinders the research in this area. Recently, a dataset, CYBY23, of Twitter threads having main tweets and the replies of Bystanders was published on Kaggle in Oct 2023. The dataset has extracted features related to toxicity and sensitivity of the main tweets and reply tweets. The authors have got manual annotators to assign the labels of Bystanders' roles. Manually labeling bystanders' roles is a labor-intensive task which eventually raises the need to have an automatic labeling technique for identifying the Bystander role. In this work, we aim to suggest a machine-learning model with high efficiency for the automatic labeling of Bystanders. Initially, the dataset was re-sampled using SMOTE to make it a balanced dataset. Next, we experimented with 12 models using various feature engineering techniques. Best features were selected for further experimentation by removing highly correlated and less relevant features. The models were evaluated on the metrics of accuracy, precision, recall, and F1 score. We found that the Random Forest Classifier (RFC) model with a certain set of features is the highest scorer among all 12 models. The RFC model was further tested against various splits of training and test sets. The highest results were achieved using a training set of 85% and a test set of 15%, having 78.83% accuracy, 81.79% precision, 74.83% recall, and 79.45% F1 score. Automatic labeling proposed in this work, will help in scaling the dataset which will be useful for further studies related to cyberbullying.

*Keywords*—*Bystanders; cyberbullying; machine learning; defender; instigator; impartial; toxicity; twitter*

## I. INTRODUCTION

With the emergence of technology in this digital era the dynamics of human connection have changed. Social media platforms have evolved into incredible tools for connecting individuals from all over the world. However, some individuals use it positively while others engage in terrible conduct on social media. The destructive phenomenon of cyberbullying has emerged as a result of the rise of social media platforms [1]. As our lives grow more entwined with the virtual domain, the frequency and consequences of cyberbullying have caught the interest of scholars, educators, and lawmakers.

Bullying is defined as a recurring pattern of hostile or aggressive behavior carried out by an individual or group that meets three criteria: repetition, intent to harm, and lack of

authority [2]. The major actors engaged in bullying irrespective of the circumstances in which it occurs are the perpetrator (bully), the victim, and bystanders. Bystanders in the cyberbullying landscape might be considered passive witnesses, which may involve strangers, who are often lured into the online chaos. They have the potential to either perpetuate or mitigate the trauma of victims. Bystanders have the potential to make a positive impact in bullying situations. Victims feel less worried and disappointed when they are surrounded by compassionate peers. Bystanders are present during bullying occurrences 80% of the time, and when they react, the bullying stops in 57% of cases within 10 seconds.

Statistics highlight a harsh reality, emphasizing the importance of acknowledging and addressing cyberbullying. According to recent surveys, an enormous percentage of people of different ages have been victims of internet abuse. Moreover, the findings provide a comprehensive picture, emphasizing the frequency of cyberbullying. Many studies use Twitter as one of the most popular data sources to identify cyberbullying as it is the most popular social networking site where cyberbullying is prevalent because of its constant conversation atmosphere which allows users to openly express their emotions, thoughts, and opinions [3].

Children and teenagers are more familiar with the internet nowadays than ever before, at younger ages. This pattern has given rise to a major concern of cyberbullying [4]. Cyberbullying has a significant impact on victims both physically and psychologically. Bullying can cause depression, anxiety, loneliness, dejection, low self-esteem, anger, self-harming behavior, alcohol and drug usage, and engagement in violence or crime. Physical health suffers as well, resulting in headaches, sleeplessness, abdominal pain, food disorders, and nausea. Cyberbullying has also shown long-term effects on victims, causing stress, continuous misery, sleep difficulties, and even issues like hunger [5].

## II. BACKGROUND AND RELATED WORK

To identify bullying, an annotation technique [6] was created to recognize textual aspects of cyberbullying, which includes posts by bullies and responses from victims and the audience. The fundamental goal of [6] research is to acquire an understanding of the language aspects of cyberbullying. This is accomplished in two stages by gathering and annotating a dataset. A harmfulness score is calculated for each post in the first phase to determine whether it is part of a cyberbullying incident. If that's the case, annotators divide the authors' roles into four categories: harasser, victim, bystander defender, and bystander assistant. A binary classifier for each fine-grained

bullying category has been built by the end. Additional features like semantic information were not explored in this research.

The study discovered that the spread of hatred from the primary posts to the replies significantly impacts how annotators identify a thread, frequently leading to reclassification as bullying rather than plain aggression [7][8]. An examination of the entire thread assists annotators in understanding the intent behind the use of specific phrases, which may have different interpretations depending on the context [9]. This finding is consistent with earlier research emphasizing the impact of bystander behavior in online environments. Bystanders' reactions are socially influenced and can be formed by their interactions with offensive comments, resulting in peer pressure and antisocial conduct. The study emphasizes the complex dynamics of online interactions, namely the involvement of bystanders in contributing to the overall classification of content as bullying. The study discovered that the spread of hatred from the primary posts to the replies significantly impacts how annotators identify a thread, frequently leading to reclassification as bullying rather than plain aggression. [7][8] An examination of the entire thread assists annotators in understanding the intent behind the use of specific phrases, which may have different interpretations depending on the context [9]. This finding is consistent with earlier research emphasizing the impact of bystander behavior in online environments. Bystanders' reactions are socially influenced and can be formed by their interactions with offensive comments, resulting in peer pressure and antisocial conduct. The study emphasizes the complex dynamics of online interactions, namely the involvement of bystanders in contributing to the overall classification of content as bullying [7][8].

The work done by [10] focuses on two objectives one is to detect cyberbullying as a binary classification problem and to detect participant roles as a multi-class classification problem. In simple terms, the focus is on evaluating the performance of models that could classify whether the post is cyberbullying-related and if it is the prediction of author's role is done. But there is a need for a more comprehensive and integrated approach that goes beyond individual posts to capture the dynamics of entire discussions in the context of cyberbullying.

While [11] contains two cyberbullying corpora in Dutch and English language. Both are manually annotated with bullying types and participant roles: harasser/bully - the individual who initiates the harassment, Victim - the one who is harassed, Bystander-Assistant: someone who assists the harasser. Bystander-defender:a person who supports the victim. This dataset has a serious problem of imbalance in the data. As "Bystander-Assistant" was the minority class, so the "Bystander-Assistant" was merged with the "Harasser" class to reduce the skew. However, there was still a large amount of imbalance between the "Harasser", "Victim" and "Defender" classes, and between "Bullying" and "No Bullying" in both English and Dutch Corpus which could negatively affect the machine learning corpus. Table II summarizes the related work in this area.

As concluded, there are many datasets available in the field of cyberbullying research on Twitter. Previous studies on cyberbullying detection as mentioned in Table I on Twitter relied on datasets labeled based on individual tweets, failing to capture the complexities of cyberbullying incidents. Labeling

the roles of bystanders is a time-consuming job, especially when examining Twitter threads with a significant number of replies, as it demands a thread-by-thread approach thereby creating a need to automate the labeling techniques.

The uniqueness of the dataset [12], [13], [14] used in this research is the inclusion of labels for bystanders' roles and aggressiveness level of Cyberbullying. Many of the existing datasets solely focus on labeling the main post lacking information about the participants involved such as Bystanders. To the best of our knowledge, this dataset is different from the existing datasets. It contains 112 Twitter threads including the main post and the replies on that post totalling around 639 tweets. It also includes the primary tweets and bystander replies. These threads are grouped by conversation ID. By incorporating efficient machine learning models on this dataset better classification can be done leading to a deeper understanding of real-world scenarios [13], [14].

Through the Literature Survey, it can be said that there are not many Twitter datasets available where bystander roles in Cyberbullying are classified. The dataset used here [12], [13], [14] contains multiple types of Bystander roles such as defender, instigator, impartial, or other. It also consists of multi-class labels either as bullying with high aggression, bullying with low aggression, or aggression without indication of bullying.

The rest of the paper is organized as follows: Section II-A presents the motivation and objectives of the proposed work. Section III explains the methodology of the research. Experiments with results and their analysis are discussed in the Section IV followed by conclusions and suggestions for future work in Section V.

TABLE I. PUBLICLY AVAILABLE DATASETS FOR CYBERBULLYING

| Data Source | Data size | Data Language | Data Gathering Tools |
|---|---|---|---|
| ASKfm[6] | 91,370 Dutch posts | Dutch | GNU Wget software |
| ASKfm[10] | - | English | AMICA |
| Facebook[15] | 100 comments | English | |
| ASKfm[4] [11] | 113,698 English, 78,387 Dutch | English and Dutch | GNU Wget software |
| Twitter[16] | 79,799 conversations with 528,041 tweets | English | Twarc |

*A. Motivation*

The risk of cyberbullying is increasing year by year due to increased access to technology, low-cost internet connections, and the leaders enthusiastically pursuing and pushing the dream of "Digital India," making its assessment and prevention even more crucial. The vast majority of people now have access to the Internet. The children and teenagers are the most susceptible members, as they are driven into cyberspace before they are psychologically capable of making sense of it. According to Microsoft's Global Youth Online Behaviour Survey, India ranks third in cyberbullying, with 53% of respondents, primarily youngsters, admitting to have experienced online bullying, trailing only China and Singapore.

TABLE II. CYBERBULLYING DETECTION, AND BULLYING TYPES

| Characteristics | Preprocessing steps | Classifier | Technique | Classification |
|---|---|---|---|---|
| Bag of words, polarity based on sentiment lexicon features [6] | Tokenization, lemmatization and PoS-tagging | Binary Classifier | SVM | Harasser, victim and bystander |
| An ensemble model is extended with a pre-trained BERT embedding layer, hidden neural layer, and a softmax output layer [10] | Replacing slang words, abbreviations, decoding emoticons, punctuations removal, upper to lower case, tokenization and special token additions | Binary Classifier | Ensemble model | Harasser, Victim, Bystander defender, Bystander assistant |
| Latent Semantic Analysis, multitask multimodality Gated Recurrent Unit, and Dirichlet Multinomial Mixture are applied to detect cyberbullying [15] | Tokenization, lemmatization, stemming, removing special characters and stop words, | Random Forest | Latent semantic analysis and feature extraction | Denigration, Trickery, Flaming, and Cyberstalking |
| Discovering bystander effect from the negative correlation between the number of Twitter users in the conversation before a toxic tweet was sent and the number of users who responded to the toxic tweet in a non-toxic manner. [16] | tweets with only links, images, and videos were discarded | - | Multivariate regression analysis, Poisson regression model, linear regression model | Bystanders |
| Multiclass classification to determine cyberbullying with Participant role detection. Investigating feature-engineered single and ensemble classifier setups and transformer-based pre-trained language models (PLMs) [11] | Tokenization, lemmatization and part-of-speech-tagging | Linear classification, Voting classifier, Cascading classifier | SVM, Logistic regression, passive-aggressive, SGD Random BL, Majority BL | Harasser, Victim, Bystander defender, Bystander assistant |

Bystanders play an important role in dealing with cyberbullying situations where they can change the dynamics of relationships. They can respond in three ways: by replicating the perpetrator's toxic behavior, by interfering with the toxic talk and sticking up for the victim, or by just observing the unfolding events. However, the mechanisms of bystander behavior in cyberspace in response to hate speech are complex. This complication emerges because the existence of other internet users may reduce one's sense of obligation to interfere, expecting that someone else will do so. Bystanders in smaller groups, on the other hand, feel a larger need to intervene in cases of cyberbullying [17].

Most of the datasets that are available publicly do not emphasize any information related to the Bystander roles in Cyberbullying. Considering the effect of the bystanders, it is important to classify its role. The motive is to explore and potentially implement automatic labeling techniques for the dataset CYBY23 [12]. The integration of automated labeling techniques in the dataset CYBY23 [12] helps to enhance the dataset's scalability and usability for future studies in cyberbullying research. The overarching goal is to contribute to the advancement of research in the field, offering insights that can foster a healthier online environment.

### B. Objective

In this work, we aim to suggest a highly efficient technique for

1) Automated labeling of bystander roles in cyberbullying tweets.
2) Finding out the most effective features extracted from the text of the tweets.

For the above objectives, we will deploy several machine learning models and experiment with various pre-processing, and feature selection techniques to discover the most efficient one among those.

### III. METHODOLOGY AND PROPOSED MODEL

In this section, the methodology of our research work is described. Flow chart for the same is given in Fig. 1. The Major steps are listed below:

1) Data Ingestion: The dataset, CYBY23, was downloaded from the Kaggle website [13], [12].
2) Data Pre-processing: Initially, the imbalance of the data was removed by using the SMOTE technique [18]. Further, data was pre-processed to make it suitable for machine learning models. The features of the main tweet were augmented with those of reply tweets and some unwanted features were removed. Categorical features were converted to numeric values.
3) Deployment of Machine Learning Models: Twelve machine-learning models [19] were deployed on the pre-processed data of Bystanders. The parameters of all the models were hypertuned to give their best performance. Pycaret library of Python [1] was used for this purpose. The models were evaluated based on accuracy, precision, recall, and F1 score metrics.
4) Experiments with Feature Selection: Next, various combinations of feature sets were experimented with like Toxicity features only (extracted from Perspective API [2]), Sensitivity features only (extracted from TextBlob [3]), and combinations of these features. Further, highly correlated features and less relevant features were removed to judge the performance of machine learning models.

Finally, a machine learning model having best accuracy and F1 score was recommended for automatic labeling of Bystanders role. The automation of Bystanders role detection will help in the early detection of cyberbullying cases and reduce their number to a greater extent.

Each of the steps involved in the process is explained below in detail:

---

[1]https://pycaret.org
[2]https://perspectiveapi.com/
[3]https://textblob.readthedocs.io/en/dev/

Fig. 1. Flow chart of methodology

### A. Dataset Description

The dataset related to bystanders was downloaded from Kaggle [12]. Alfurayj et al. [13] used Twitter API to extract 1024 tweets from January 2022 to January 2023. 150 tweet threads were collected. Information such as the date of the tweet, tweet ID, screen name of the user and user ID associated with the tweet, number of likes & retweets, and text of the tweet was downloaded. Religion, ethnicity, sarcasm, and racial orientation were among the keywords and hashtags used to crawl this information, which could lead to harassment remarks. A manual annotation process for the labeling of Bystanders was used. Annotators followed the guidelines given in [20] and assessed the aggressiveness of individual tweets, identified bystander roles in replies, and made higher-level judgments about the overall aggressiveness of the thread after considering the main post, replies, and bystander roles. Following the annotation process, threads lacking agreement from at least five annotators were eliminated, reducing the tweets to 639. The dataset, meeting the criteria for a good dataset, contained a minimum of 10% to 20% bullying cases, with cyberbullying with high aggression representing only 11.6%. Instigators were notably high in both bullying categories. The investigation focused on bystander contagion risk, with a higher prevalence of instigators associated with instances of bullying, as evidenced by the dataset. They realized the need for the automation of annotation for labeling of Bystanders' role because of the labor-intensive nature of manual annotation and hence a dataset, named CYBY23, was uploaded on the Kaggle website [12] for public use. CYBY23 dataset had the Twitter threads containing both the main posts and the replies from Bystanders. Each tweet had the text of the tweet along with certain general features of

the tweets. Further, they extracted the Toxicity features using Perspective API and sentiment features using TextBlob for each tweet. There were 639 tweets in the dataset with the labels of bystanders' roles (manually annotated).

So, the dataset, CYBY23 [12], had six general features, namely, tweet_id, reply_id, text , created_at, favorite_count, retweet_count for each tweet. Six features were derived from Perspective API , namely, Insult, Threat, Identity_Attack , Profanity ,Toxicity , and Severe_Toxicity, and three features were derived from TextBlob , namely, polarity, subjectivity, and sentiment. Feature 'class label' was assigned to the main tweet only and the feature 'bystander role label' was assigned to the reply tweet only. Thus, the dataset had sixteen features for main tweets and fifteen features for reply tweets. (see Table III).

TABLE III. FEATURES OF ORIGINAL DATASET CYBY23

| General | Perspective API | TextBlob | Main Tweet | Reply Tweet |
|---|---|---|---|---|
| tweet_id | Insult | polarity | class label | bystander role label |
| reply_id | Threat | subjectivity | | |
| text | Identity _Attack | sentiment | | |
| created _at | Profanity | | | |
| favorite _count | Toxicity | | | |
| retweet _count | Severe _Toxicity | | | |

TABLE IV. FEATURES OF PRE-PROCESSED DATASET

| General | Perspective API (Main Tweet) | Perspective API (Reply Tweet) | TextBlob (Main Tweet) | TextBlob (Reply Tweet) | Main Tweet | Reply Tweet |
|---|---|---|---|---|---|---|
| favorite _count | Insult _main | Insult | Polarity _main | Polarity | Class label | Bystander role label |
| favorite _count _main | Threat _main | Threat | subjectivity _main | Subjectivity | | |
| retweet _count | Identity _Attack _main | Identity _Attack | Sentiment _main | sentiment | | |
| retweet _count _main | Profanity _main | Profanity | | | | |
| | Toxicity _main | Toxicity | | | | |
| | Severe _Toxicity _main | Severe _Toxicity | | | | |

### B. Data Preprocessing

Certain pre-processing steps were applied to the CYBY23 dataset [12] before running the machine-learning models. Those are listed below:

1) The feature 'bystander role label' had four string values, namely, "This person agrees with the main post (instigator)", "This person disagrees with the main post (defender)", "This person is not taking any sides (impartial)" and "This person posted unrelated

replies (Other)". These string values were converted to numeric values between 0 to 3.

2) To study the effect of the main tweet on the reply tweets, the features of the main tweet were concatenated with the features of the reply tweet, and a new dataset was created. The new dataset had seven general features, six toxicity-related features of the reply tweet and main tweet, three sentiment-related features of the reply tweet and main tweet, feature 'class label' of the main tweet, and feature 'bystander role label' of the reply tweet. Thus, the new dataset had 28 features. Names of main tweet features were suffixed with _main. Since main tweets were concatenated column-wise with reply tweet, so the number of total tweets reduced from 639 to 524.

3) The features tweet_id, reply_id, and created_at were removed as they were not required for the models. So new dataset had 25 features for all the tweets.

4) The feature 'text' was removed from the dataset, because toxicity features using Perspective API and sentiments features using TextBlob had already been computed from the 'text' feature. Thus, the new dataset had 24 features for all the tweets.

After pre-processing, we got the dataset having 524 tweets and 24 features for each tweet (see Table IV). Out of these 24 features, the feature 'bystander role label' was used as the target feature for all machine learning models.

### C. Model Development

In this work, we deployed different machine learning models [19] using Pycaret library. A brief description of each of the models is given below:

- AdaBoost Classifier (ADA): Adaptive Boosting Classifier is an ensemble classifier, that benefits from training several weak classifiers and then combining the result, with more weightage given to the classifier that gives more accuracy.

- Decision Tree Classifier (DT): A flowchart-like tree structure where each internal node denotes a test on an attribute, each branch represents the outcome of a test, and each leaf node holds a class label.

- Extra Trees Classifier (et): An ensemble machine learning method based on decision trees. The dataset sampling for each tree is done randomly, without replacement. The features subset is also assigned randomly to each tree.

- Gradient Boosting Classifier (GBC): This classifier is an additive model of decision trees and is often employed for both regression and classification tasks.

- K Neighbors Classifier(KNN): A learning method that uses the nearest neighbors to classify a data point.

- Linear Discriminant Analysis (LDA): A method used to find a linear combination of features that best separates two or more classes in a dataset.

- Light Gradient Boosting Machine (LGBM) & Extreme Gradient Boosting (EGB): Both are gradient boosting frameworks that use tree-based learning algorithms. They are recognized for their efficiency and predictive accuracy.

- Logistic Regression (LR): A foundational statistical method to model the probability of a certain class or event based on one or multiple predictor features.

- Naive Bayes (NB): A probabilistic classifier based on applying Bayes' theorem, it assumes independence between features.

- Random Forest Classifier(rf): An ensemble learning method that uses decision trees. Each decision tree comprises of dataset drawn by bootstrap sampling. The 'majority voting' is used to make final prediction.

- Ridge Classifier (RC): A classification algorithm that employs L2 regularization. It can help prevent overfitting and often delivers better performance in scenarios with multicollinearity.

- SVM - Linear Kernel(SVM): A learning method that finds a hyperplane to separate the two classes such that it maximizes predictive accuracy while avoiding over-fitting.

### D. Model Validation

The proposed model was validated using various feature selection techniques:

1) Experimenting on various types of features (Toxicity Based, Sentiment Based)
2) Removal of Highly correlated features
3) Removal of less significant features
4) Hypertuning the parameters of machine learning models

Model efficiency was analyzed after applying each of the techniques mentioned above.

### E. Model Evaluation

1) Use of the Tool: We used Pycaret Python library which speeds up the process of experiments related to machine learning and empowers us to run multiple ML models simultaneously. It also helps in hypertuning the parameters of the models which gives us the best performance.
2) Evaluation Metrics: Four metrics, accuracy, precision, recall, and F1 score are used to evaluate the models. Using a wide range of evaluation metrics caters to various aspects of prediction quality.
3) Cross-Validation: We applied K-Fold cross-validation. This method partitioned the training data into 'K' subsets, training on 'K-1' of them and validating on the remaining subset. This process was iteratively executed until each subset had been used for validation, offering a robust average performance metric.
4) Various Train-Test Split: Various splits for training and test sets were used to validate the model.

## IV. Experimental Results and Analysis

This section presents the experimental setup, their results, and analysis.

### A. Platforms Used

We used Python using Jupyter Notebook and Google collaboratory for running the experiments. Pycaret library was used to run the machine learning models. Plotting of graphs was done using Matplotlib and Pandas library.

### B. Dataset

A pre-processed dataset (see Table IV), having 524 tweets and 24 features for each tweet, was used in further experiments.

*1) Handling Imbalance of Dataset:* The class distribution of the dataset having 524 tweets is shown in Fig. 2 (a). High imbalance can be observed in the number of instances of unique values of the target feature 'bystander role label'. Imbalance can be handled by undersampling or oversampling the minority class. However, undersampling has the chance of losing important information. So, we used an oversampling technique, Synthetic Minority Oversampling Technique (SMOTE) [18] to handle the imbalance. SMOTE generates synthetic samples for the minority class and creates a balanced dataset. Fig. 2 (b) depicts the balanced dataset with 912 data points.



Fig. 2. (a) Class distribution (b) Class distribution after resampling (SMOTE).

### C. Model Deployment using Various Feature Selection Techniques

We experimented with different feature selection techniques on various machine learning models. Pycaret was used to run all the models. The models were evaluated using accuracy, precision, recall, and F1 score metrics. The results of running all machine learning models using Pycaret are shown in Table V. The experiments and their results are mentioned below:

- Case 1: Initially we run the experiments using only the toxicity features derived from Perspective API. Random Forest Classifier(rf), Gradient Boosting Classifier (gbc), Light Gradient Boosting Machine (lightgbm), and Extra Trees Classifier (et) performed best, each with accuracy as well as F1 score of 72% (approx).

- Case 2: Further, the experiments were run on Sentiments features derived from TextBlob. Approximately 70% accuracy, and 70% F1 score were achieved using Random Forest Classifier(rf), Gradient Boosting Classifier (gbc), Light Gradient Boosting Machine (lightgbm) and Extra Trees Classifier (et) classifier (see Table V).

- Case 3: Next, we experimented with both the toxicity features (mentioned in case 1) and sentiments features (mentioned in case 2). With this feature set, accuracy as well as F1 score of approx. 75% was achieved with all the four classifiers mentioned in Case 1 and Case 2. Thus, indicating that instead of using only Toxicity or Sentiment features, results are better when both are used.

- Case 4: From the feature set mentioned in case 3, we computed the correlation coefficient among features (see Fig. 3). We found that the feature Severe_Toxicity_main is highly correlated to Toxicity_main. Similarly, the features Profanity and Toxicity, favorite_count_main and retweetcount_main, Toxicity and Insult, Severe_Toxicity and Profanity, Toxicity_main and Insult_main are highly correlated. Thus, we removed the features, 'Severe_Toxicity_main', 'Profanity', 'Toxicity', 'favorite_count_main', 'Toxicity_main, and were left with 19 features.
  After removing the correlated features, the highest accuracy of 76% was achieved. Again, the same four classifiers, rf, gbc, lightgbm, and et, performed best.



Fig. 3. Heatmap showing correlation among features.

- Case 5: Next, we experimented with finding the importance of the features mentioned in case 3. Some feature ClassLabel_main, sentiment, sentiment_main, and retweet_count were ruled out (see Fig. 4) because of their low importance.
  After removing the less important feature, we checked the efficiency of our models (see Table V). Random Forest Classifier(rf) performed best with 76% accuracy and 78% F1 score.

- Case 6: Further, we chose a feature set that was formed after removing the highly correlated features as well

Fig. 4. Feature importance.



Fig. 6. Comparison of accuracy for different feature sets.

as the less important features from the features given in Case 3. Running all the machine learning models using Pycaret gave the results mentioned in Table V. We observe that Random Forest Classifier(rf) again performed best with 77.6% accuracy and 79.8% F1 score.

Fig. 5 compares the accuracy of all the classifier models for each of the feature set case discussed above.



Fig. 5. Comparison of accuracy for different models.

For each of the feature set case discussed above, Fig. 6 compares the accuracy achieved by the different classifier models. Here, results from only those classifier are plotted, which achieved more than 50% accuracy.

As is discussed above for all the six cases and is evident from Fig. 5 and Fig. 6, Random Forest Classifier(rf) performs best for most of cases. Also, the best result is achieved by the feature set formed by including both the Toxicity and Sentiments features and by removing both the least significant features and the highly correlated features.

### D. Feature Importance

Before going for further experimentation, we would like to give some observations related to the importance of features as depicted in Fig. 4.

1) We found that the features ClassLabel_main, sentiment, sentiment_main, and retweet_count have very less importance as compared to other features (see Fig. 4). This indicates that the level of aggression of the whole thread denoted by ClassLabel_main has little impact on the model performance. The sentiment of the reply tweet and the sentiment of the main tweet has very little role to play along with the number of retweets indicated by retweet_count.

2) Features Insult and Toxicity have the highest importance. One of them can be considered an important feature since they are highly correlated.

3) Feature Threat of the main tweet and reply tweet is almost equally important.

4) Features Identity_Attack, Profanity, Insult, Toxicity, Severe_Toxicity, Polarity, Sentiment of the main tweet have low importance as compared to the corresponding features of the reply tweet except for the Threat and Subjectivity feature.

5) Comparing the set of features based on Perspective API and TextBlob[4], we can observe that features based on Perspective API have more importance.

Summarizing the observations, the top features among all are Toxicity, Identity_Attack, Threat_main, Profanity, and Threat.

### E. Different Train-test Split

Then we experimented with different train-test splits for judging the performance of Random Forest classifier. The data

---

[4]https://textblob.readthedocs.io/en/dev/

TABLE V. EVALUATION METRICS FOR DIFFERENT COMBINATIONS

| Model | Accuracy | Recall | Precision | F1-Score |
|---|---|---|---|---|
| **Ada Boost Classifier** | | | | |
| Toxicity Features Only | 0.4797 | 0.4797 | 0.5374 | 0.4875 |
| Sentiment Features Only | 0.4608 | 0.5408 | 0.5489 | 0.5292 |
| Toxicity & Sentiment Features | 0.4766 | 0.4766 | 0.5430 | 0.4838 |
| Removing Correlated Features | 0.4872 | 0.5172 | 0.56 | 0.5222 |
| Removing Less Relevant Features | 0.4964 | 0.4764 | 0.5243 | 0.4791 |
| Removing Correlated and Less Relevant | 0.5031 | 0.5031 | 0.5606 | 0.51 |
| **Decision Tree Classifier** | | | | |
| Toxicity Features Only | 0.6644 | 0.6444 | 0.6435 | 0.6392 |
| Sentiment Features Only | 0.6600 | 0.6600 | 0.6609 | 0.6566 |
| Toxicity & Sentiment Features | 0.6740 | 0.6740 | 0.6746 | 0.6693 |
| Removing Correlated Features | 0.6756 | 0.6756 | 0.6763 | 0.6725 |
| Removing Less Relevant Features | 0.6944 | 0.6944 | 0.6988 | 0.6909 |
| Removing Correlated and Less Relevant | 0.6982 | 0.6912 | 0.7011 | 0.6885 |
| **Extra Trees Classifier** | | | | |
| Toxicity Features Only | 0.7045 | 0.7445 | 0.7476 | 0.7436 |
| Sentiment Features Only | 0.6954 | 0.7054 | 0.7150 | 0.7012 |
| Toxicity & Sentiment Features | 0.7186 | 0.7634 | 0.7662 | 0.7611 |
| Removing Correlated Features | 0.7217 | 0.7617 | 0.7665 | 0.759 |
| Removing Less Relevant Features | 0.7266 | 0.7666 | 0.7713 | 0.7641 |
| Removing Correlated and Less Relevant | 0.7418 | 0.7618 | 0.7668 | 0.7605 |
| **Gradient Boosting Classifier** | | | | |
| Toxicity Features Only | 0.7290 | 0.7290 | 0.7383 | 0.7287 |
| Sentiment Features Only | 0.7178 | 0.7178 | 0.7236 | 0.7155 |
| Toxicity & Sentiment Features | 0.7315 | 0.7571 | 0.7590 | 0.7552 |
| Removing Correlated Features | 0.7373 | 0.7273 | 0.7278 | 0.7236 |
| Removing Less Relevant Features | 0.7415 | 0.7415 | 0.7445 | 0.7401 |
| Removing Correlated and Less Relevant | 0.7563 | 0.7163 | 0.7246 | 0.714 |
| **K Neighbors Classifier** | | | | |
| Toxicity Features Only | 0.6008 | 0.6708 | 0.6976 | 0.6722 |
| Sentiment Features Only | 0.5856 | 0.6456 | 0.6644 | 0.6414 |
| Toxicity & Sentiment Features | 0.6597 | 0.6597 | 0.6765 | 0.6596 |
| Removing Correlated Features | 0.6633 | 0.6333 | 0.6393 | 0.6282 |
| Removing Less Relevant Features | 0.6648 | 0.6348 | 0.6516 | 0.636 |
| Removing Correlated and Less Relevant | 0.6706 | 0.6206 | 0.613 | 0.6128 |
| **Linear Discriminant Analysis** | | | | |
| Toxicity Features Only | 0.4984 | 0.4984 | 0.5117 | 0.4926 |
| Sentiment Features Only | 0.4862 | 0.5062 | 0.4799 | 0.4773 |
| Toxicity & Sentiment Features | 0.5301 | 0.5501 | 0.5560 | 0.5435 |
| Removing Correlated Features | 0.5376 | 0.5376 | 0.5439 | 0.5329 |
| Removing Less Relevant Features | 0.5556 | 0.5156 | 0.5255 | 0.5127 |
| Removing Correlated and Less Relevant | 0.5687 | 0.4987 | 0.517 | 0.4991 |
| **Light Gradient Boosting Machine** | | | | |
| Toxicity Features Only | 0.7136 | 0.7336 | 0.7358 | 0.7311 |
| Sentiment Features Only | 0.6931 | 0.7131 | 0.7187 | 0.7090 |
| Toxicity & Sentiment Features | 0.7233 | 0.7633 | 0.7648 | 0.7617 |
| Removing Correlated Features | 0.7337 | 0.7337 | 0.7368 | 0.73 |
| Removing Less Relevant Features | 0.7495 | 0.7695 | 0.7726 | 0.7676 |
| Removing Correlated and Less Relevant | 0.7587 | 0.7587 | 0.7678 | 0.757 |
| **Logistic Regression** | | | | |
| Toxicity Features Only | 0.4404 | 0.4404 | 0.4339 | 0.4219 |
| Sentiment Features Only | 0.4953 | 0.4953 | 0.4662 | 0.4593 |
| Toxicity & Sentiment Features | 0.5000 | 0.5000 | 0.4883 | 0.4815 |
| Removing Correlated Features | 0.5579 | 0.5579 | 0.5389 | 0.5409 |
| Removing Less Relevant Features | 0.5601 | 0.4701 | 0.4741 | 0.4666 |
| Removing Correlated and Less Relevant | 0.5722 | 0.5222 | 0.52 | 0.5099 |
| **Naive Bayes** | | | | |
| Toxicity Features Only | 0.4907 | 0.4907 | 0.5382 | 0.4473 |
| Sentiment Features Only | 0.4984 | 0.4984 | 0.5520 | 0.4436 |
| Toxicity & Sentiment Features | 0.5485 | 0.5485 | 0.5881 | 0.5115 |
| Removing Correlated Features | 0.5593 | 0.5393 | 0.5775 | 0.4991 |
| Removing Less Relevant Features | 0.5691 | 0.4591 | 0.471 | 0.4204 |
| Removing Correlated and Less Relevant | 0.5791 | 0.4891 | 0.5066 | 0.4399 |
| **Random Forest Classifier** | | | | |
| Toxicity Features Only | 0.7226 | 0.7226 | 0.7233 | 0.7200 |
| Sentiment Features Only | 0.7085 | 0.7085 | 0.7111 | 0.7036 |
| Toxicity & Sentiment Features | 0.7346 | 0.7346 | 0.7514 | 0.7442 |
| Removing Correlated Features | 0.7462 | 0.7462 | 0.7773 | 0.7661 |
| Removing Less Relevant Features | 0.7606 | 0.7606 | 0.7869 | 0.7778 |
| Removing Correlated and Less Relevant | 0.7766 | 0.7666 | 0.7987 | 0.7938 |
| **Ridge Classifier** | | | | |
| Toxicity Features Only | 0.4890 | 0.4890 | 0.4739 | 0.4686 |
| Sentiment Features Only | 0.5016 | 0.5016 | 0.4623 | 0.4531 |
| Toxicity & Sentiment Features | 0.5439 | 0.5439 | 0.5182 | 0.5155 |
| Removing Correlated Features | 0.5455 | 0.5455 | 0.5264 | 0.521 |
| Removing Less Relevant Features | 0.5547 | 0.5047 | 0.4969 | 0.4889 |
| Removing Correlated and Less Relevant | 0.5697 | 0.5097 | 0.5067 | 0.4952 |
| **SVM Linear Kernel** | | | | |
| Toxicity Features Only | 0.2397 | 0.2397 | 0.1374 | 0.1305 |
| Sentiment Features Only | 0.2321 | 0.2321 | 0.1438 | 0.1190 |
| Toxicity & Sentiment Features | 0.2492 | 0.2492 | 0.2285 | 0.1450 |
| Removing Correlated Features | 0.2541 | 0.2541 | 0.1733 | 0.1568 |
| Removing Less Relevant Features | 0.2664 | 0.2664 | 0.2388 | 0.1434 |
| Removing Correlated and Less Relevant Features | 0.2765 | 0.262 | 0.2297 | 0.1971 |

is the case with F1 score. Hence, best accuracy of 78.83% and F1 score of 79.45% is reported to be achieved at 85% training set and 15% test set (see Table VI).



Fig. 7. Model performance on different train-test split.

TABLE VI. MODEL PERFORMANCE FOR VARIOUS TRAIN-TEST SPLITS

| Train-Test Split | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| 60-40 | 0.7150 | 0.7241 | 0.7050 | 0.7119 |
| 65-35 | 0.7281 | 0.7509 | 0.7081 | 0.7277 |
| 70-30 | 0.7189 | 0.7506 | 0.7010 | 0.7204 |
| 75-25 | 0.7280 | 0.7637 | 0.7080 | 0.7308 |
| 80-20 | 0.7486 | 0.7846 | 0.7286 | 0.7520 |
| 85-15 | 0.7883 | 0.8179 | 0.7483 | 0.7945 |
| 90-10 | 0.7883 | 0.8126 | 0.7434 | 0.7945 |

## V. CONCLUSION AND FUTURE WORK

In this paper, a machine learning model for automatic labeling of Bystanders detection has been proposed. Initially, Pycaret was used to find the best model using the features mentioned in the CYBY23 dataset [12]. Later, various feature selection techniques have been used to increase the efficiency of the model. The proposed model has been validated by using different train-test splits. The results of various combinations has been discussed in length. Finally, the Random Forest classifier with a training set of 85% and 15% has been chosen as the best model for Bystanders detection. Further, the Importance of various features from the given dataset has been discussed. Despite the best efforts of applying machine learning techniques for the given dataset CYBY23 [12], the authors feel that the small size of the dataset hinders the research in this area. The achieved results will be more promising on a larger dataset.

The research work in the future may be directed toward increasing the dataset size and finding a more efficient model for automatic labeling. The dataset size can be increased by extending the work to other social media posts. Various ways of finding the sentiments can be used using Natural Language Processing techniques. Deep learning models can be experimented with for the deployment of an efficient model for automatic labeling. The mentioned dataset can be regarded as a multi-label dataset with two class labels, namely aggression level, and bystanders role and further experiments can be performed in that direction.

was split in multiple split percentages, starting from 60% till 95% with a window of 5%. Fig. 7 summarizes the results of running the Random forest classifier when the training set is split from 60% till 95% with a window of 5%. We observe that the accuracy increases with the increase in the size of the training set but almost stabilizes when it reaches 85%. Similar

REFERENCES

[1] T. Mahlangu, C. Tu, and O. Pius, "A review of automated detection methods for cyberbullying," in *International Conference on Intelligent and Innovative Computing Applications (ICONIC)*. IEEE, 12 2018, pp. 1–5.

[2] H. Kallmen and M. Hallgren, "Bullying at school and mental health problems among adolescents: a repeated cross-sectional study." *Child Adolesc Psychiatry Ment Health*, vol. 74, 2021. [Online]. Available: https://doi.org/10.1186/s13034-021-00425-y

[3] S. Salawu, Y. He, and J. Lumsden, "Approaches to automated detection of cyberbullying: A survey," *IEEE Transactions on Affective Computing*, vol. 11, no. 1, pp. 3–24, Mar. 2020.

[4] C. Van Hee, G. Jacobs, C. Emmery, B. Desmet, E. Lefever, B. Verhoeven, G. De Pauw, W. Daelemans, and V. Hoste, "Automatic detection of cyberbullying in social media text," *PLOS ONE*, vol. 13, no. 10, p. e0203794, Oct. 2018.

[5] P. K. Smith, J. Mahdavi, M. Carvalho, S. Fisher, S. Russell, and N. Tippett, "Cyberbullying: its nature and impact in secondary school pupils," *Journal of Child Psychology and Psychiatry*, vol. 49, no. 4, pp. 376–385, 2008. [Online]. Available: https://acamh.onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-7610.2007.01846.x

[6] C. Van Hee, E. Lefever, B. Verhoeven, J. Mennes, B. Desmet, G. De Pauw, W. Daelemans, and V. Hoste, "Automatic detection and prevention of cyberbullying," in *International Conference on Human and Social Analytics, Proceedings*, P. Lorenz and C. Bourret, Eds. IARIA, 10 2015, pp. 13–18.

[7] M. Tsvetkova and M. Macy, "The social contagion of antisocial behavior," *Sociological Science*, vol. 2, pp. 36–49, 02 2015.

[8] E. J. Villota and S. G. Yoo, "An experiment of influences of facebook posts in other users," in *2018 International Conference on eDemocracy and eGovernment (ICEDEG)*, Ambato, Ecuador, 2018, pp. 83–88.

[9] K. Yokotani and M. Takano, "Social contagion of cyberbullying via online perpetrator and victim networks," *Computers in Human Behavior*, vol. 119, p. 106719, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0747563221000418

[10] G. Rathnayake, T. Atapattu, M. Herath, G. Zhang, and K. Falkner, "Enhancing the identification of cyberbullying through participant roles," in *Proceedings of the Fourth Workshop on Online Abuse and Harms*. Association for Computational Linguistics, 11 2020, pp. 89–94. [Online]. Available: https://aclanthology.org/2020.alw-1.11

[11] Jacobs, Gilles and Van Hee, Cynthia and Hoste, Veronique, "Automatic classification of participant roles in cyberbullying: can we detect victims, bullies, and bystanders in social media text?" *Natural Language Engineering*, vol. 28, no. 2, pp. 141–166, 2022. [Online]. Available: http://doi.org/10.1017/s135132492000056X

[12] H. Alfurayj, N. S. Yee, and S. L. Lutfi, "Cyberbullying bystander dataset 2023," 2023. [Online]. Available: https://www.kaggle.com/dsv/6486152

[13] H. S. Alfurayj, S. L. Lutfi, and N. S. Yee, "Bystanders unveiled: Introducing a comprehensive cyberbullying corpus with bystander information," *TENCON 2023 - 2023 IEEE Region 10 Conference (TENCON)*, pp. 1012–1017, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:265354251

[14] H. S. Alfurayj and S. L. Lutfi, "Exploring bystanders' roles in labeled cyberbullying threads on twitter: A preliminary analysis," *TENCON 2023 - 2023 IEEE Region 10 Conference (TENCON)*, 2023.

[15] R. C. J.I. Sheeba, S. Pradeep Devaneyan, "Identification and classification of cyberbully incidents using bystander intervention model," *International Journal of Recent Technology and Engineering*, vol. 8, no. 2S4, p. 1–6, Aug. 2019. [Online]. Available: http://dx.doi.org/10.35940/ijrte.B1001.0782S419

[16] A. Aleksandric, M. Singhal, A. Groggel, and S. Nilizadeh, "Understanding the bystander effect on toxic twitter conversations," *ArXiv*, vol. abs/2211.10764, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:253734314

[17] M. Obermaier, N. Fawzi, and T. Koch, "Bystanding or standing by? how the number of bystanders affects the intention to intervene in cyberbullying," *New Media & Society*, vol. 18, no. 8, pp. 1491–1507, 2016.

[18] N. Chawla, K. Bowyer, L. Hall, and W. Kegelmeyer, "Smote: Synthetic minority over-sampling technique," *J. Artif. Intell. Res. (JAIR)*, vol. 16, pp. 321–357, 06 2002.

[19] I. Sarker, "Machine learning: Algorithms, real-world applications and research directions," *SN Computer Science*, vol. 2, 03 2021.

[20] Van Hee, Cynthia and Verhoeven, Ben and Lefever, Els and De Pauw, Guy and Daelemans, Walter and Hoste, Veronique, "Guidelines for the fine-grained analysis of cyberbullying, version 1.0," 2015.

# Using Deep Learning to Recognize Fake Faces

Jaffar Atwan[1], Mohammad Wedyan[2], Dheeb Albashish[3], Elaf Aljaafrah[4], Ryan Alturki[5], Bandar Alshawi[6]

Prince Abdullah bin Ghazi Faculty of Information and Communication Technology
Al-Balqa Applied University, Jordan[1,3,4]
Department of Computer Sciences-Faculty of Information Technology and Computer Science
Yarmouk University, Irbid, 21163, Jordan[2]
Department of Software Engineering-College of Computing
Umm Al-Qura University, Makkah, Saudi Arabia[5]
Department of Computer and Network Engineering-College of Computing
Umm Al-Qura University, Makkah, Saudi Arabia[6]

*Abstract*—In recent times, many fake faces have been created using deep learning and machine learning. Most fake faces made with deep learning are referred to as "deepfake photos." Our study's primary goal is to propose a useful framework for recognizing deep-fake photos using deep learning and transformative learning techniques. This paper proposed convolutional neural network (CNN) models based on deep transfer learning methodologies in which the designed classifier using global average pooling (GAP), dropout, and a dense layer with two neurons that use SoftMax are substituted for the final fully connected layer in the pretrained models. DenseNet201, the suggested framework, produced the best accuracy of 86.85% for both the deepfake and real picture datasets, while MobileNet produced a lower accuracy of 82.78%. The obtained experimental results showed that the proposed method outperformed other state-of-the-art fake picture discriminators in terms of performance. The proposed architecture helps cybersecurity specialists fight deepfake-related cybercrimes.

*Keywords—Deep learning; machine learning; deepfake; convolutional neural network; global average pooling*

## I. INTRODUCTION

Artificial intelligence (AI) is the process of creating devices that mimic human intelligence in terms of behaviour and thought. The term can also refer to any device exhibiting characteristics of the human mind, like problem-solving and learning [1]. An ideal attribute of AI is the capacity to simplify and carry out actions that are most likely to achieve a specific goal. A subset of AI is machine learning (ML). Massive volumes of unstructured data, including text, photos, and videos, are ingested by deep-learning algorithms to allow this autonomous learning. ML aims to replicate how humans learn and increases in accuracy over time by using data and algorithms [2, 3].

ML is a crucial component of the developing field of data science. Algorithms are trained in data-mining projects to categorize, forecast, and unearth important insights using statistical techniques. With the goal of influencing important growth metrics, these insights drive decisions within applications and organizations [4]. As big data continues to grow and improve, data scientists will become increasingly in demand. It should be possible to use ML to find the information needed to answer many important business questions. Deep learning can be classified as a subset of machine learning [5]. Deep learning uses less complex concepts than those employed in

ML and uses artificial neural networks that are designed to imitate human brain networks. Previously, the intricacy of neural networks has been limited by computer power. Larger and more complex neural networks are now conceivable due to advancements in big data analytics, which allow computers to see, learn from, and respond to complex events more quickly than people can. Deep learning makes it possible to categorize images, identify faces, translate languages, recognize audio, and determine whether a face is real or fake. It can tackle pattern recognition issues and does not require human intervention [6, 7].

The face is a person's most recognizable feature. The security hazards of facial modification are becoming increasingly more significant because of the rapid development of facial synthesis technology. Several algorithms based on deep-learning techniques can replace one person's face with another person's realistic-looking visage [8]. Additionally, new AI technology called deepfake combines the faces of two different people. A number of methods based on generative adversarial networks (GANs) produce high-resolution deepfake images that are more accurate than previous technologies [9]. This is cause for concern, as deepfake information can circulate quickly due to the rise of mobile phones and the emergence of multiple social networking sites [10]. Initially, deepfake photos could be distinguished by the human eye because of a pixel collapse phenomenon that tends to produce unnatural visual contrasts in skin tones and face features. However, over time and with the development of technology, deepfakes have essentially merged with natural imagery [11].

Deepfake techniques frequently require enormous volumes of audio, video, or image data to produce convincing photographs that look natural. However, while deepfakes represent huge development in technological capability, there are some negatives. There is a prevalence of deepfakes of public people, including athletes, politicians, and celebrities, in the abundance of films and photos that can be found online [12]. Additionally, deepfake technologies can be used to ridicule and humiliate people. Deepfakes are considered to be the most harmful sort of synthetic media. They utilize celebrities' voices and photos without their permission to make political or humorous content about them. Due to the simplicity of the numerous applications making deepfakes, anyone can use this technology to make artificial content that is indistinguishable from actual content. It is not only public people who can be affected by deepfake

technology. One use of deepfake content is cyberbullying, which affects a large population of young people [10]. A number of factors are taken into account in the sophisticated approach of deepfake image detection. The basic steps of image classification include identifying a suitable classification scheme, collecting training patterns, image pre-processing, feature extraction, choosing a suitable evaluation method, and evaluating accuracy.

The remainder of the essay is structured as follows. Section II provides background on deepfakes, GANs, and a summary of a range of studies and previous research on image classification. Section III focuses on research procedures and methods of work. It includes a detailed explanation of the models used. Section IV presents an experimental setup. Section V describes the results of the experiment obtained using the selected dataset on a set of models and makes a comparison between them based on several criteria. Finally, Section VI presents conclusions and suggestions for further work.

## II. Literature Review

The development of technology has made life easier in many respects. However, there have also been instances where technology has been abused, which has resulted in some serious issues. One example is digital image technology. There are many tools and software that make it is easy to modify any digital image. For example, anyone with even a basic understanding of Photoshop can quickly and simply create a fake image of another person [13].

There has been a lot of recent research on the use of these kinds of forgeries. Advancements in the disciplines of AI means that people may now alter a raw image and use it in both positive and harmful ways because, crucially, these techniques can provide incredibly life-like outcomes. This introduced us to the realm of deepfake pictures [14]. For example, [15] uses deep learning as a technology that creates face recognition and can determine whether a profile image is authentic or not, with the aim of finding a reliable method to distinguish between actual and phony. This study included real and fake face detection utilizing deep learning methods built on neural networks in two image datasets. They chose the ResNet50 model as the best match and used a trained dataset of 9,000 photos. The training accuracy was 99.18%. The research in [16], transfer learning methods from previously trained depth models like ResNet50 and VGG16 were used in the proposed model and three benchmark datasets were used to assess the proposed model. The findings collected demonstrate that the suggested model outperforms the current models. The study in [17] used enhanced datasets for real and fake face identification to compare the most popular modern face-recognition classifiers, including Custom convolutional neural network (CNN), VGG19, and DenseNet-121. They found performance can be increased while using fewer computational resources due to data augmentation. According to the authors preliminary findings, VGG19 outperforms all other examined models and has a maximum accuracy of 95%. To create ensemble-like multi-attention networks for detecting deep fake media, this work attempts to provide a complete examination of the mentioned methods, structures, and mechanisms.

The research in [18] attempts to address the difficulty of differentiating between real and fake pictures by developing an algorithm that can distinguish between real and fake pictures. The algorithm used in [18] seeks to differentiate between real images and deep fakes. The dataset was tested against five transfer learning methods as well as an 18-layered bespoke CNN model that was described in the research. The proposed model was able to test with an accuracy of 98.77%, whereas InceptionV3 produced the best results of the transfer learning models with a testing accuracy of 97.10%. Comparing deep-fake and real photos, the unique CNN model performed better than any other model previously employed. The main goal of [19] was to develop a reliable and accurate method for recognizing deepfake images. The significance of this work lies in obtaining positive outcomes while utilizing the CNN architecture. This study employed eight CNN architectures to identify deep-fake images from big datasets. The findings were accurate and dependable. For some criteria, like F1 score, precision, and area under the Receiver operating characteristic ROC curve, the custom model used in this investigation performed marginally better than VGG Face in terms of recall.

The research in [20] provided a pipeline for categorizing and recognizing human faces from input visual samples. The second stage employed a number of deep learning (DL)-based techniques to calculate deep features from the returned faces. A support vector machine (SVM), a type of classifier, was trained on these characteristics to assess whether the data was real or fake. They compared the performance of numerous feature extractors based on their published results and found that DenseNet169 and its SVM classifier surpassed the competition. Table I summarizes the previously mentioned studies.

## III. Materials and Methods

In order to detect fake faces, this work builds a group of pre-trained models with fine-tuning. A final choice is made for a testing image by fine-tuning five pre-trained models (DenseNet201, MobileNet, InceptionV3, ResNet50, and Xception) and fusing their projected probabilities. The pre-trained models use transfer learning to reduce their weights so that they can perform a similar classification problem. For the classification of faces, ensemble learning of previously trained models achieves greater results.

### A. Pretrained Dense Net

A variation on the ResNet design is the densely linked convolutional network (DenseNet) architecture suggested by [21]. In this architecture, layers are connected to one another using the summation technique. In comparison to the ResNet design, the summing operation aids in further improving generalization ability and better resolving the issue of the vanishing gradient. The features that are taken from each layer are used as input for the following layers in this method. Reusing feature maps could help the overall performance be improved even further. The architecture of DenseNet201 contains 201 layers, hence the name. In this paradigm, high performance can be attained with little memory and little computational expense. DenseNet comes in a variety of sizes, including 121, 169, 201, and 264.

### B. Pretrained MobileNet

The Google research team created the MobileNet architecture [22] for object identification on portable devices.

TABLE I. Summary of the Most Important Classification Studies on Fake Faces

| Authors | Dataset used | CNN architectures |
|---|---|---|
| Maher Salman et al. [15] | Real and fake faces detection | VGG16, ResNet50 InceptionV3, MobileNet |
| Taeb et al. [17] | Real and fake face detection 140K real and fake faces | VGG19, DenseNet121 |
| Sharma et al. [16] | 140k real and fake faces Fakefaces Real and fake face detection | VGG16, ResNet50 |
| Dhar [18] | 140K real and fake faces | VGG16, DenseNet121 InceptionV3, VGG19, ResNet50 |
| Shad et al. [19] | 140K real and fake faces | DenseNet201, DenseNet169, ResNet50, VGG16, VGG19, VGGFace |
| Masood et al. [20] | DeepFake Detection Challenge (DFDC) | VGG16, VGG19, ResNet101, Inception V3, DenseNet-169, InceptionResV2, XceptionNet, MobileNetv2, EfficientNet, NASNetMobile |

MobileNet architecture presented a depth-wise separable convolution along with 11 point-wise convolution layers, having 32 times fewer parameters compared to conventional convolutions. MobileNet architecture outperformed VGG16 achieving higher accuracy during training on ImageNet dataset and requiring 27 times less computational power. Through depth-wise convolution, one depth-wise kernel was employed all the input channel. Point-wise convolution utilizes 11-bit kernel size CONV layer to calculate a linear combination of several input channels. The preceding method reduces the feature maps dimensionality significantly.

### C. Pretrained Resnet

The ResNet50 network has a lot of depth. With it, more complicated networks can be constructed (which might refer to as networks inside networks) utilizing common network components known as residual modules and train them using stochastic gradient descent (SGD). The ResNet architecture [23] was groundbreaking work that demonstrated how residual modules can be used to train very deep networks using regular SGD. By applying identity-mapping techniques to update the residual coefficients, accuracy can be attained. Its architecture drastically reduces its size by using a global average pooling layer rather than a fully linked layer. This network is called ResNet50 because the architecture has 50 levels.

### D. Pretrained Xception

Xception architecture, which stands for extreme inception and was introduced by François Chollet [24], is an improvement on the Inception design. In this architecture, the initialization modules from the Inception design are replaced by residual connections and depth-wise separable convolutions. It is possible that the depth-wise separable convolution will lower memory and processing expenses. The Xception architecture consists of 14 modules, each with 36 convolutional layers. All connections between modules, except for the first and last, are created via linear residual connections.

### E. Pretrained Inception

The third iteration of the Inception model, the Inception V3 architecture [24] has a total of 159 layers. Instead of utilizing a single kernel size (such as 3x3 or 5x5), the Inception module uses several convolution sizes, such as 1x1, 3x3, and 5x5 filter sizes. The fundamental concept behind using various convolution sizes is that it enables the extraction of multi-level characteristics from the input image during each convolution process. Pointwise 11 convolution is also employed in this architecture to cut down on the number of parameters. The computational cost is decreased by the pointwise convolutions. The network has undergone numerous iterations due to its ongoing evolution. InceptionV1, InceptionV2, InceptionV3, InceptionV4, and Inception-Resent are common variants. Table II shows the summary of the deep architectures employed in this study.

### F. Experimental Design

The proposed method for identifying fake or real faces based on the CNN architecture is described in this section. By using five different models, this study attempts to create a deep-learning model for face classification. The entire workflow of suggested solution is depicted in Fig. 1. The diagram illustrates the three basic steps of the model. The first phase is loading the dataset and image processing, the second is using the pre-trained model to extract features, and the third is using the selected features and classifying images. The proposed model uses datasets as input, and the final output is to classify images and evaluate and visualize the results.

Five different deep learning models – ResNet50, Inception V3, DenseNet201, Xception, and MobileNet – have been used as the base models and pre-trained for classification using the ImageNet dataset. An approach called transfer learning is used to train these models. In transfer learning, a pre-trained network performs better than a network that was trained from scratch. As shown, constructing classification solutions with transfer learning is quicker and more effective than doing it without. CNN also plays a fascinating role in classification. Two components make up each model: a feature extractor and a classifier. The classifier is used to categorize the collected features, whereas the feature extractor works to extract features using a convolutional base layer. In order to determine if the output is a fake face or a real face, The convolutional base layers and adapt the final classification layer are kept by adding new sets of layers such as global average pooling (GAP), dropout, and the dense layer.

### IV. Experimental Setup

### A. Datasets

The proposed model on a deepfake and real images dataset acquired from the Kaggle website is tested. https://www.kaggle.com/datasets/manjilkarki/deepfake-and-real-images. Five CNN models were trained to distinguish between fake and real images. The dataset is divided into a training set and a test set. The training set has 4,700 images, of which 2,500 are real, and the rest are fake. The testing set has 540 images, of which

TABLE II. A SUMMARY OF THE DEEP ARCHITECTURES EMPLOYED IN THIS STUDY

| Architecture | Convolutional layer count | Count of face centred cubic (FCC) layers | Parameter count for training |
|---|---|---|---|
| DenseNet201 | 199 | 2 | 20.2 million |
| MobileNet | 53 | 3 | 3.4 million |
| ResNet50 | 48 | 2 | 25.6 million |
| Xception | 70 | 1 | 22.9 million |
| InceptionV3 | 42 | 1 | 23.9 million |



Fig. 1. Proposed experimental design.

300 are real and the rest are fake. Both real and altered photos can be found in this dataset. The faces, which are produced using a variety of techniques, are modified images. To extract the most value from these photos, this dataset was processed. Each picture is a 256x256 jpg image of a real or fake human face.

### B. Pre-processing

The most important part of the model is the pre-processing method. To minimize overfitting, data augmentation was implemented. A 224x224x3 image was provided as the last input for the recommended model. Image augmentation is the process of creating new training samples from existing ones. To create a new pattern, slightly modify the original image; for example, you can make the new image slightly brighter or crop part of the original image. The original image can be mirrored to produce a new one [25]. There are various techniques to help increase the number of data points, such as rotation, shift, zooming, and horizontal fling. Augmented datasets were used for these experiments. To make the expanded dataset better fit the trained models, and they were scaled it and added horizontal flips by added a shifting of 0.1, a zoom range of 0.5, and a 45-degree rotation to the datasets.

### C. Extraction of Features

In the feature extraction approach, the network of convolutional and pooling layers that serve as the extraction of features were kept while removing the fully connected layers of a pretrained CNN model. The feature extractor can be expanded with fully linked layers and machine-learning classifiers. As a result of the dataset being more appropriate for this model, the network's performance on it is improved. Also, the final fully connected layer and retrieved features were kept with the trained models ResNet50, Inception V3, DenseNet201, Xception, and MobileNet.

### D. Classification

Deep features were extracted and sent through the ResNet50, Inception V3, DenseNet201, Xception, and MobileNet models before being transferred to user-specific layers that were specifically designed for them. Deep features that had been concatenated were scaled in one-dimension (1D) form using GAP, producing feature maps that were appropriate for the succeeding two completely connected layers. Two fully connected layers and introduced dropout (0.5) in the midst of the fully connected layers were used to improve efficiency and generalize learning. The activation function and the output are ultimately produced by a dense layer with two neurons that uses the SoftMax activation function for binary classification.

### E. Evaluation Criteria

In this study, the TensorFlow package, Keras API, and Python programming were used to implement all the pre-trained models (DenseNet201, MobileNet, ResNet50, Xception, and Inception V3). Additionally, Google Colab Pro was used for all tests. The model is trained and optimized using the Adam optimizer. A cycle of updating network weights using all the training data is known as an epoch. A model's performance will advance over time as the number of epochs rises. All models were tested across 25 epochs with a learning rate of 0.001 and a batch size of 32. Dropout was introduced to expedite training, enhance learning, boost precision, and avoid overfitting. The inputs used to train the model are shown in the Table III.

1) Accuracy: The percentage of correctly categorized images is what is meant by accuracy. TP + TN / (TP + TN + FP + FN).

2) Precision: It is the proportion of positively anticipated categories to positively classed categories that were effectively recognized. TP/ (TP + FP).

TABLE III. HYPERPARAMETERS USED IN THE SUGGESTED TRANSFER LEARNING MODELS

| Hyperparameters | Value |
|---|---|
| Image size | 224 x 224 |
| Optimizer | Adam |
| Learning rate | 0.001 |
| Batch size | 32 |
| Dropout | 0.5 |
| Number of epochs | 25 |
| Activation function | SoftMax |

3) Recall: The recall rate is the proportion of subjects who were correctly classified out of all positively classified subjects. TP / (FN + TP).

4) F1 score: The F1 score is typically employed to make it possible to measure both precision and recall simultaneously. The harmonic mean is used in place of the arithmetic mean. As a result, the penalize extreme values more. 2*(precision*recall)/ (precision + recall)

## V. THE RESULTS

The proposed methods were used to test out a number of pre-trained deep-learning models that were available. Performance in various models was enhanced using different optimizers. A number of CNN models (see Table IV) were implemented using the deepfake and real images dataset. This demonstrated good facial image classification accuracy. Additionally, the figures of each proposed model were shown and explained using the dataset. The automatic identification and classification of faces are presented in depth in this part, along with the results of the studies. To create a reliable classifier, numerous trials were carried out with list models, InceptionV3, DenseNet201, MobileNet, ResNet50, and Xception. This study's primary objective is to evaluate the effectiveness of deep learning architectures. On the basis of performance metrics for precision, recall, and F1 score, the five designs employed in the study were assessed. The experimental results attained for each model are shown in Table IV. The results table show that the classifier performs well for each class.

Table IV and Fig. 2 show the results for the accuracy, precision, and F1-score recall of the deep fake and real images dataset, which includes two classes of fake faces and real faces using five of the pre-training models with optimizer Adam, anumber of epochs of 25 for each model with a SoftMax activation function, and a batch size of 32. The model that achieved the highest accuracy was DenseNet201, with a rate of 86.58% and the highest recall of 0.86, a precision of 0.87, and an F1 score of 0.87, while ResNet50 had an accuracy of 83.33%. The accuracy for Xception was 84.07% and, 85.0% for Inception V3. The MobileNet model provided relatively low accuracy, sensitivity, precision, and F1-score values for all classes.

Graph (A) from Fig. 2 shows the accuracy of the model DenseNet201 throughout training and validation over a period of 25 epochs. As the number of epochs rises, the accuracy of training and validation appears to increase. However, there are some variations in the validation accuracy over time. The validation accuracy fell below 65% in the first three epochs.

However, the results approached a score of 86% by the 25th epoch, while the validation loss fluctuated, eventually falling to zero across the remaining epochs.

Fig. 3 (A) shows the accuracy of the model MobileNet throughout training and validation over a period of 25 epochs. As the number of epochs rises, the accuracy of training and validation appears to increase. However, there are some variations in the validation accuracy over time. The validation accuracy fell below 66% in the first 15 epochs. However, the results approached a score of 82% by the 25th epoch, while the validation loss fluctuated, eventually falling to zero across the remaining epochs.

Fig. 4 graph (A) shows the accuracy of the model ResNet50 throughout training and validation over a period of 25 epochs. As the number of epochs rises, the accuracy of training and validation appears to increase. However, there are some variations in the validation accuracy over time. The validation accuracy fell below 55 in the first five epochs. However, the results approached a score of 83% score by the 25th epoch, while the validation loss fluctuated, eventually falling to zero across the remaining epochs.

Fig. 5 graph (A) shows the accuracy of the model Xception throughout training and validation over a period of 25 epochs. As the number of epochs rises, the accuracy of training and validation appears to increase. However, there are some variations in the validation accuracy over time. The validation accuracy fell below 72% in the first 10 epochs. However, the results approached a score of 84% by the 25th epoch, while the validation loss fluctuated, eventually falling to zero across the remaining epochs.

Fig. 6 graph (A) shows the accuracy of the model InceptionV3 throughout training and validation over a period of 25 epochs. As the number of epochs rises, the accuracy of training and validation appears to increase. However, there are some variations in the validation accuracy over time. The validation accuracy fell below 68% in the first 15 epochs. However, the results approached a score of 85% by the 25th epoch, while the validation loss fluctuated, eventually falling to zero across the remaining epochs.

### A. Performance Evaluation Metrics

There is a concept known as a confusion matrix in the context of machine learning, deep learning, and, more specifically, the issue of statistical classification. A table that summarizes how well a classification model works on a collection of test data or real values from the set is known as a confusion matrix. A result, the algorithm's performance can be assessed and commonalities between classes can be quickly found. In further detail, the confusion matrix is a clear account of the outcomes of a categorization task that contains a summary of the right and wrong predictions. The true negative (TN) condition occurs when the model predicts the negative class accurately. The negative type in this instance relates to an actual face. A false negative (FN) occurs when the model forecasts the negative class inaccurately and incorrectly predicts that the face was real. A false positive (FP) occurs when the model forecasts the positive class inaccurately; that is, it predicted the face to be a fake but it was incorrect. When the model accurately predicts the positive class, it is said to be a true

TABLE IV. THE EXPERIMENTAL RESULTS OBTAINED ON THE DEEPFAKE AND REAL IMAGES DATASET USING MODELS

| Pretrained models | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| DenseNet201 | 86.58% | 0.87 | 0.86 | 0.87 |
| MobileNet | 82.78% | 0.83 | 0.83 | 0.83 |
| ResNet50 | 83.33% | 0.83 | 0.83 | 0.83 |
| Xception | 84.07% | 0.85 | 0.83 | 0.84 |
| InceptionV3 | 85.00% | 0.87 | 0.84 | 0.84 |



Fig. 2. The accuracy (A) and Loss (B) of the DenseNet201 model during training and validation.



Fig. 3. The accuracy (A) and Loss (B) of the MobileNet model during training and validation.

Fig. 4. The accuracy (A) and Loss (B) of the ResNet50 model during training and validation.



Fig. 5. The accuracy (A) and Loss (B) of the Xception model during training and validation.

positive (TP). The positive category in this instance refers to a fake face.

Fig. 7 displays the confusion matrix for the DenseNet201 model for the deepfake and real images dataset. The number of images is 540, divided into 240 fake images and 280 real images. Forty-five images were incorrectly labeled as fake when they were real faces and 26 images were real but incorrectly labeled as fake. Furthermore, 195 of the photographs were accurately identified as fake, while 274 of the images were correctly identified as real.

Fig. 8 displays the confusion matrix for the MobileNet model for the deepfake and real images dataset. The number of images is 540, divided into 240 fake images and 280 real im-

ages. Forty-five images were incorrectly labeled as fake when they were real faces and 48 images were real but incorrectly labeled as fake. Furthermore, 195 of the photographs were accurately identified as fake, while 252 of the images were correctly identified as real.

Fig. 9 displays the confusion matrix for the ResNet50 model for the deepfake and real images dataset. The number of images is 540, which were divided into 240 fake images and 280 real images. Forty-seven images were incorrectly labeled as fake when they were real faces and 43 images were real but incorrectly labeled as fake. Furthermore, 193 of the photographs were accurately identified as fake, while 257 of the images were correctly identified as real.
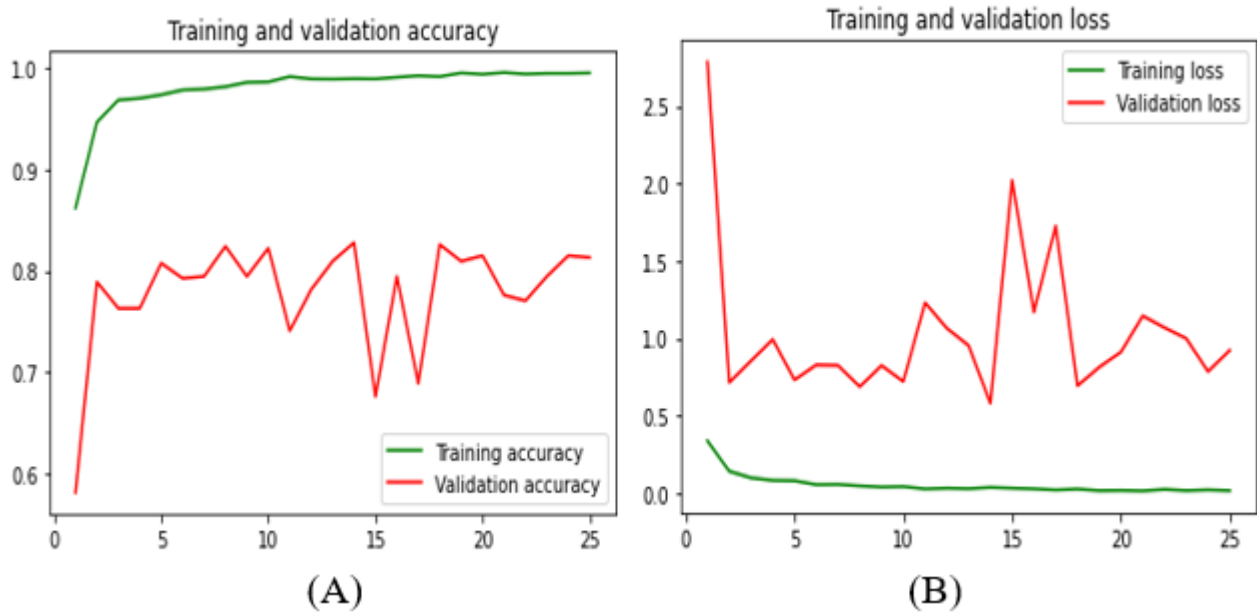
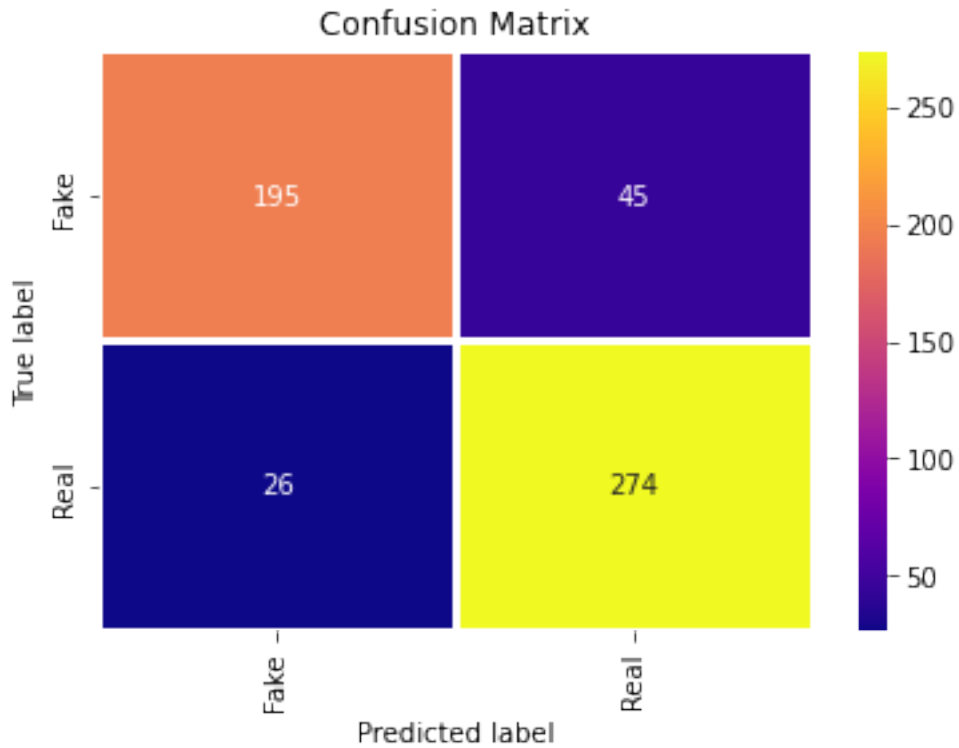Fig. 6. The accuracy (A) and Loss (B) of the InceptionV3 model during training and validation.



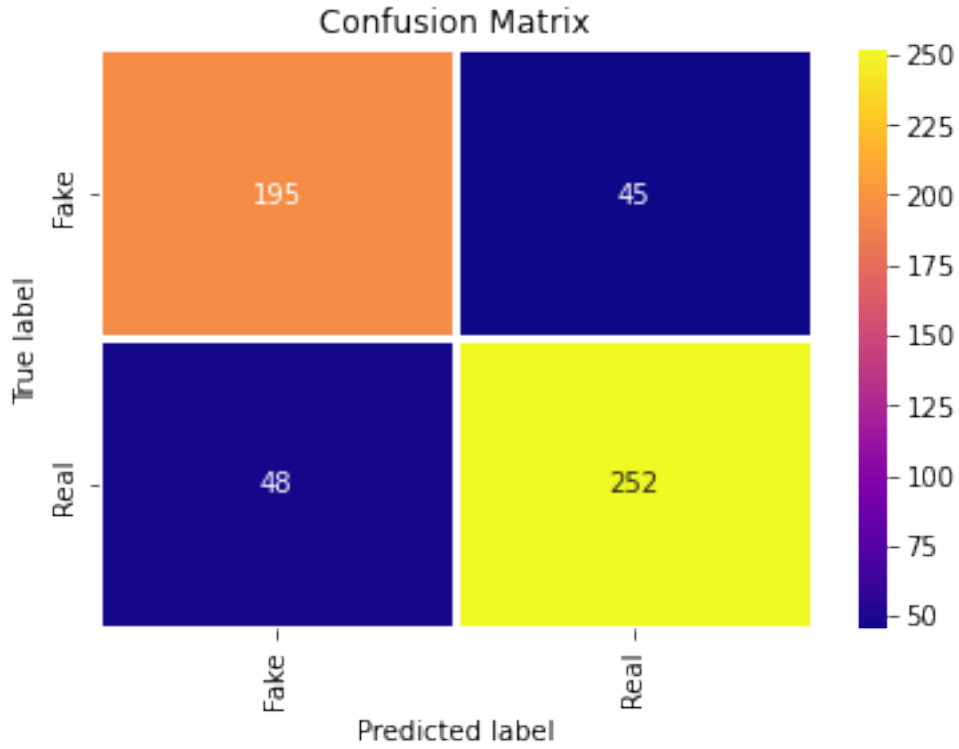Fig. 7. The result of the prediction of the DenseNet201.

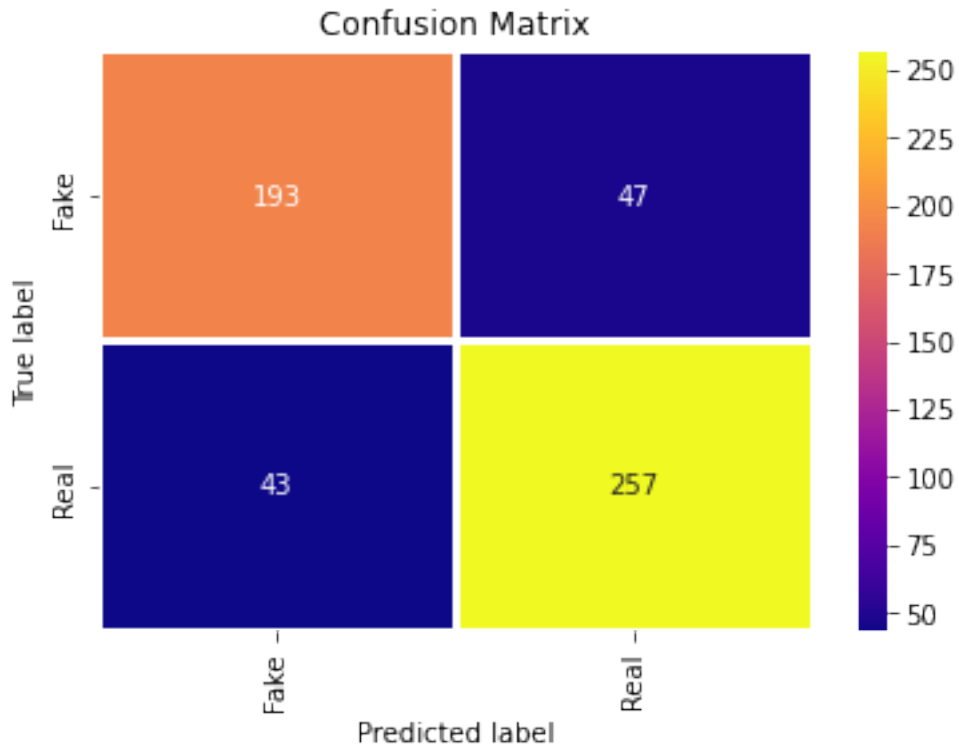Fig. 8. The result of the prediction of the MobileNet.



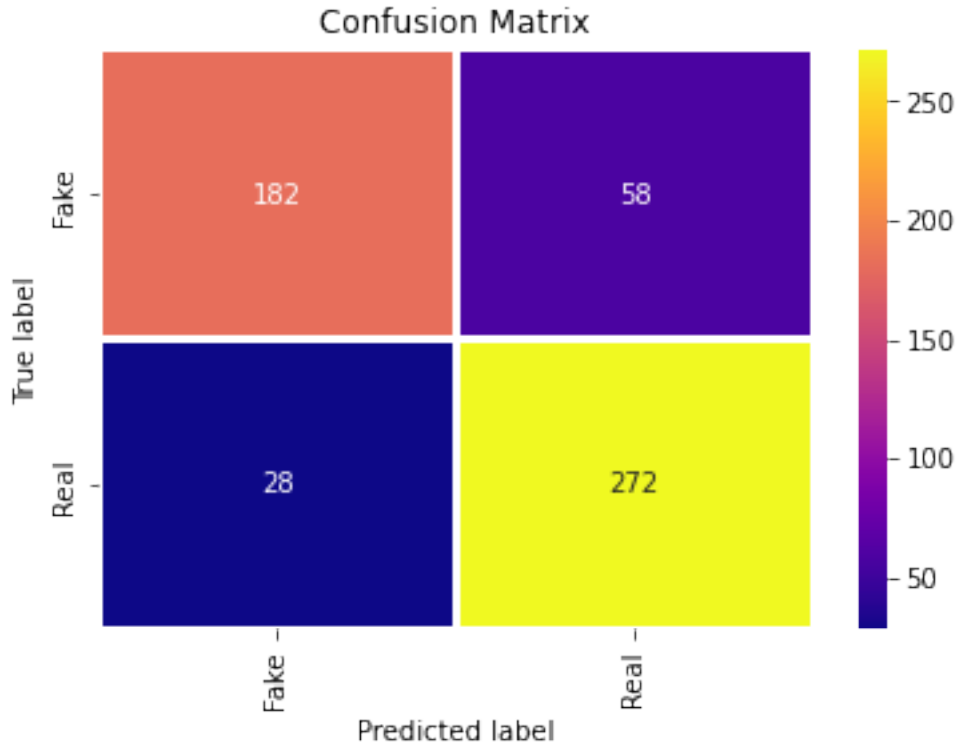Fig. 9. The result of the prediction of the ResNet50.

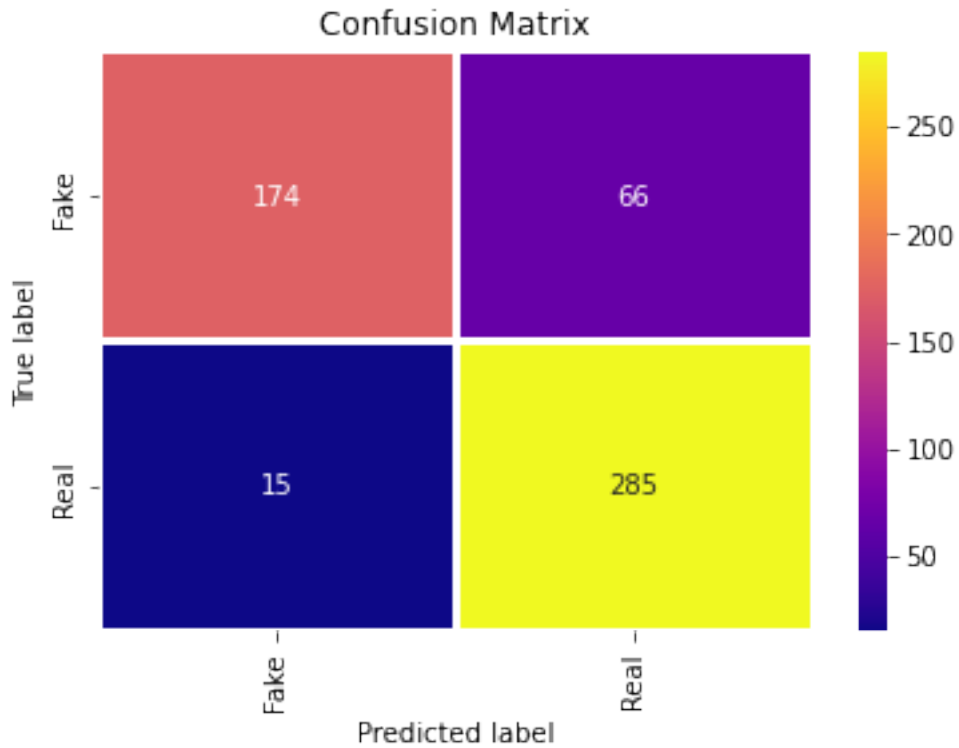Fig. 10. The result of the prediction of the Xception.



Fig. 11. The result of the prediction of the InceptionV3.

Fig. 10 displays the confusion matrix for the Xception model for the deepfake and real images dataset. The number of images is 540, which were divided into 240 fake images and 280 real images. Fifty-eight images were incorrectly labeled as fake when they were real faces and 26 images were real but incorrectly labeled as fake. Furthermore, 182 of the photographs were accurately identified as fake, while 272 of the images were correctly identified as real.

Fig. 11 displays the confusion matrix for the Inception V3 model for the deepfake and real images dataset. The number of images is 540, which were divided into 240 fake images and 280 real images. Sixty-six images were incorrectly labeled as fake when they were real faces and 16 images were real but incorrectly labeled as fake. Furthermore, 174 of the photographs were accurately identified as fake, while 285 of the images were correctly identified as real.

## VI. Conclusion and Future Works

A new technique called "deepfake" is being employed to uses AI to generate realistic but fake images of people, particularly public figures. While not all fake information is harmful, some of it genuinely threatens the global community and should be identified. The main goal of this research was to develop a reliable and accurate method for spotting phony pictures. Researchers have used a number of techniques to find deep-fake content. However, the significance of this work lies in obtaining positive outcomes while utilizing the CNN architecture. In this study, the paper employed transfer-learning techniques in the proposed framework to enhance the accuracy of detection and reduce execution time. Also, the paper applied the proposed model to several pre-trained models, and a comparison was made in terms of accuracy, sensitivity, recall, and F1 score. This research used five pre-trained models — Resnet50, Inception V3, DenseNet201, Xception, and MobileNet — to detect deep-fake images using public datasets. The dataset contained deepfake and real images, with 4,700 training images and 540 test images. The final fully connected layer in the pre-trained models was eliminated in this study and replaced with a classifier that uses dropout, GAP, and a dense layer with two neurons that employs SoftMax. Image augmentation techniques were also used, with the help of the optimizer Adam. Some improvements can be made to the deep-learning framework used in this paper, such as applying the framework to different datasets, performing experiments using pre-trained models different from those used in this paper, and merging two CNN models with each other. The aim of this research was to design an application that detects deepfakes and gives an accurate and automatic performance evaluation. Additionally, we intend to evaluate this work on low-resolution, low-light imagery and extend it to real and fake video recognition.

## References

[1] J. McCarthy, "From here to human-level AI," *Artificial Intelligence*, vol. 171, no. 18, pp. 1174–1182, 2007.

[2] J. Atwan, M. Wedyan, Q. Bsoul, A. Hamadeen, R. Alturki, and M. Ikram, "The effect of using light stemming for Arabic text classification," *International Journal of Advanced Computer Science and Applications*, vol. 12, no. 5, 2021.

[3] I. E. Naqa and M. J. Murphy, "What is machine learning?, in machine learning in radiation oncology," *machine learning in radiation oncology, Cham: Springer*, pp. 3–11, 2015.

[4] R. F. Murray, "Classification images: A review," *Journal of vision*, vol. 11, no. 5, pp. 2–2, 2011.

[5] A. Kamilaris and F. X. Prenafeta-Boldú, "Deep learning in agriculture: A survey," *Computers and electronics in agriculture*, vol. 147, pp. 70–90, 2018.

[6] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[7] M. Ibrahim, M. Wedyan, R. Alturki, M. A. Khan, and A. Al-Jumaily, "Augmentation in healthcare: Augmented biosignal using deep learning and tensor representation," *Journal of Healthcare Engineering*, vol. 2021, 2021.

[8] H. Li, Z. Lin, X. Shen, J. Brandt, and G. Hua, "A convolutional neural network cascade for face detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5325–5334.

[9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.

[10] T. T. Nguyen, Q. V. H. Nguyen, D. T. Nguyen, D. T. Nguyen, T. Huynh-The, S. Nahavandi, T. T. Nguyen, Q.-V. Pham, and C. M. Nguyen, "Deep learning for deepfakes creation and detection: A survey," *Computer Vision and Image Understanding*, vol. 223, p. 103525, 2022.

[11] T. Jung, S. Kim, and K. Kim, "DeepVision: Deepfakes Detection Using Human Eye Blinking Pattern," *IEEE Access*, vol. 8, pp. 83 144–83 154, 2020.

[12] M. Westerlund, "The emergence of deepfake technology: A review," *Technology innovation management review*, vol. 9, no. 11, 2019.

[13] B. U. Mahmud and A. Sharmin, *Deep Insights of Deepfake Technology : A Review*, vol. 5, 2020.

[14] D. Güera and E. J. Delp, "Deepfake video detection using recurrent neural networks," *15th IEEE international conference on advanced video and signal based surveillance (AVSS)*, pp. 1–6, 2018.

[15] F. M. Salman and S. S. Abu-Naser, "Classification of Real and Fake Human Faces Using Deep Learning," *International Journal of Academic Engineering Research*, vol. 6, 2022.

[16] J. Sharma, S. Sharma, V. Kumar, H. S. Hussein, and H. Alshazly, "Deepfakes Classification of Faces Using Convolutional Neural Networks," *Traitement du Signal*, vol. 39, pp. 1027–1037, 2022.

[17] M. Taeb and H. Chi, "Comparison of Deepfake Detection Techniques through Deep Learning," *Journal of Cybersecurity and Privacy*, vol. 2, pp. 89–106, 2022.

[18] A. Dhar, P. Acharjee, L. Biswas, S. Ahmed, and A. Sultana, "Detecting deepfake images using deep convolutional neural network," 2021.

[19] H. S. Shad, "Comparative Analysis of Deepfake Image Detection Method Using Convolutional Neural Network," *Comput Intell Neurosci*, vol. 2021, 2021.

[20] M. Masood, "Classification of Deepfake Videos Using Pre-trained Convolutional Neural Networks," *2021 International Conference on Digital Futures and Transformative Technologies*, 2021.

[21] G. Huang, Z. Liu, L. V. D. Maaten, and K. Q. Weinberger. [Online]. Available: https://github.com/liuzhuang13/DenseNet

[22] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.

[23] K. He, X. Zhang, S. Ren, and J. Sun. [Online]. Available: http://image-net.org/challenges/LSVRC/2015/

[24] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.

[25] X. Qiu, "Pre-trained models for natural language processing: A survey," *Sci China Technol Sci*, vol. 63, pp. 1872–1897, 2020.

# Enhancing Adversarial Defense in Neural Networks by Combining Feature Masking and Gradient Manipulation on the MNIST Dataset

Mr. Ganesh Ingle, Dr. Sanjesh Pawale
Department of Computer Engineering, Vishwakarma University, Pune, India

*Abstract*—This research investigates the escalating issue of adversarial attacks on neural networks within AI security, specifically targeting image recognition using the MNIST dataset. Our exploration centered on the potential of a combined approach incorporating feature masking and gradient manipulation to bolster adversarial defense. The main objective was to evaluate the extent to which this integrated strategy enhances network resilience against such attacks, contributing to the advancement of more robust AI systems. In our experimental framework, we utilized a conventional neural network architecture, integrating various levels of feature masking alongside established training protocols. A baseline model, devoid of feature masking, functioned as a comparative standard to gauge the efficacy of our proposed technique. We assessed the model's performance in standard scenarios as well as under Fast Gradient Sign Method (FGSM) adversarial assaults. The outcomes provided significant insights. The baseline model demonstrated a high test accuracy of $98\%$ on the MNIST dataset, yet it showed limited resistance to adversarial incursions, with accuracy diminishing to $60\%$ under FGSM onslaughts. Conversely, models incorporating feature masking exhibited a reciprocal relationship between masking proportion and accuracy, counterbalanced by an enhancement in adversarial resilience. Specifically, a $10\%$ masking ratio achieved a $96\%$ accuracy rate coupled with a $75\%$ robustness against attacks, a $30\%$ masking led to a $94\%$ accuracy with an $80\%$ robustness level, and a $50\%$ masking threshold resulted in a $92\%$ accuracy, attaining the apex of robustness at $85\%$. These results affirm the efficacy of feature masking in augmenting adversarial defense, highlighting a pivotal equilibrium between accuracy and resilience. The study lays the groundwork for further investigations into refined masking methodologies and their amalgamation with other defensive strategies, potentially broadening the scope of neural network security against adversarial threats. Our contributions are significant to the realm of AI security, showcasing an effective strategy for the development of more secure and dependable neural network frameworks.

*Keywords*—*Feature masking; neural networks; gradient manipulation; adversarial resilience; fast gradient sign method*

## I. INTRODUCTION

Enhancing adversarial defense in neural networks, particularly for image recognition tasks like those involving the MNIST dataset, can be effectively addressed by integrating feature masking and gradient manipulation. This combined approach leverages the strengths of both methods to fortify the network against adversarial attacks.

Feature Masking: This technique modifies or conceals certain features in the input data. In the context of the MNIST dataset, which comprises images of handwritten digits, feature masking could involve partially obscuring these digits. This

strategy prevents the neural network from becoming overly reliant on specific features, thus reducing its vulnerability to adversarial attacks. Research has shown that diversifying the features used by a model for classification enhances its robustness [18][19][20].

Gradient Manipulation: Neural networks adjust their parameters based on the gradient of the error relative to their current parameters. Adversarial attacks often manipulate these gradients to deceive the model. Altering the gradients, through methods like noise addition, modification, gradient clipping, or smoothing, can make the network less susceptible to minor input variations typical in adversarial attacks [5][12].

By combining feature masking and gradient manipulation, a more resilient defense against adversarial attacks can be achieved. Feature masking ensures the model does not fixate on certain input features, and gradient manipulation renders the learning process less predictable and more resistant to gradient-based adversarial methods. This holistic approach is crucial for tasks like the MNIST dataset, where inputs are relatively simple and uniform, necessitating a robust and generalizable model.

The fundamental challenge lies in the vulnerability of neural networks to adversarial perturbations. Adversarial attacks exploit the model's reliance on specific features and manipulate gradients during the learning process, leading to misclassifications. The aim of this research is to fortify neural networks against such attacks, particularly in the context of the MNIST dataset, by integrating feature masking and gradient manipulation. This paper discusses the importance of diversifying features to prevent overreliance on specific aspects of the input data and explores various gradient manipulation techniques, such as noise addition, modification, gradient clipping, or smoothing, and their potential to enhance model resilience [13][14][15][16].

Also the synergistic effects of integrating feature masking and gradient manipulation for a more comprehensive defense strategy are studies to see the impact of combining feature masking and gradient manipulation in creating a holistic defense mechanism against adversarial threats. This research aims to contribute to the development of robust neural network models, particularly for image recognition tasks like those involving the MNIST dataset. By addressing the vulnerability of neural networks to adversarial attacks through the combined approach of feature masking and gradient manipulation, the proposed methodology seeks to enhance the overall security and reliability of image recognition systems.As the field of

neural network security advances, it becomes imperative to devise comprehensive defense strategies. This paper introduces a novel approach that leverages feature masking and gradient manipulation to fortify neural networks against adversarial attacks, with a specific focus on image recognition tasks using the MNIST dataset. The research questions and objectives outlined in this paper guide the investigation into the effectiveness of this combined approach, aiming to contribute to the ongoing efforts in enhancing the security of neural networks in practical applications. The objective of the our research is to investigate the effectiveness of feature masking in preventing neural networks from fixating on specific features in the MNIST dataset. Describes the experimental setup for evaluating feature masking and its impact on model fixation,explore various gradient manipulation techniques to render the learning process less predictable and more resistant to adversarial attacks and to evaluate the combined approach of feature masking and gradient manipulation in creating a more resilient defense against adversarial attacks on neural networks trained on the MNIST dataset. Significance of the our research highlights the contribution of the proposed methodology in advancing the field of neural network security, particularly in the context of image recognition tasks. Emphasizes the potential impact on real-world applications and the broader implications for enhancing the reliability of neural network systems. This research article establishes itself as a cornerstone in advancing neural network security, presenting a holistic and innovative approach that transcends the immediate context of the MNIST dataset. The integrated feature masking and gradient manipulation methodology stands as a transformative blueprint for enhancing the security and reliability of neural network systems, with broad applications across diverse domains.

## II. BACKGROUND AND MOTIVATION

Evolution of Neural Networks: Neural networks have evolved remarkably over the last few decades, becoming more complex and powerful. They are particularly adept at image recognition tasks, outperforming traditional algorithms in most benchmarks.Rise of Adversarial Attacks: With the growing reliance on neural networks, their susceptibility to adversarial attacks has become evident. An adversarial attack involves subtly altering the input data (like images) in a way that leads the network to make incorrect predictions or classifications, while the changes remain imperceptible to the human eye.The MNIST dataset, comprising hand-written digits, is a foundational benchmark in the field of machine learning for image recognition tasks[21][35]. The simplicity and uniformity of this dataset make it an ideal testbed for studying neural network behaviors, including their vulnerability to adversarial attacks.The primary motivation is to ensure the security and reliability of neural networks in critical applications. In contexts like medical diagnosis, autonomous driving, or facial recognition, the consequences of erroneous decisions due to adversarial attacks can be severe.Improving adversarial defense helps in understanding the limitations and weaknesses of current neural network models. This understanding is crucial for developing more robust and generalizable AI systems.Enhancing adversarial defense aligns with the broader goals of AI safety and ethics. It ensures that AI systems perform reliably and safely, even in the presence of poten-

tially malicious inputs.Addressing the challenge of adversarial attacks inspires new research directions in neural network architecture design, training methodologies, and general AI robustness.As AI becomes more pervasive, regulatory bodies are increasingly focusing on the robustness and security of AI systems. Enhancing adversarial defense is thus also motivated by the need to comply with emerging regulations and standards in AI governance. The drive to enhance adversarial defense in neural networks is fueled by the need for secure, reliable, and ethical AI systems, particularly in applications where the stakes are high. The MNIST dataset serves as a fundamental platform for testing and developing these enhancements due to its simplicity and widespread use in the AI community.

## III. RELATED WORK

In recent times, the security of machine learning models has been increasingly threatened by a phenomenon known as adversarial attacks. These attacks cleverly manipulate the models by introducing subtle, often undetectable alterations, known as "adversarial examples". These alterations are designed to mislead the models into making erroneous predictions. In response to this critical issue, the scientific community has been proactive in devising a range of defensive strategies to mitigate the risks posed by these attacks.

### A. Machine Learning Security

In the field of machine learning security, recent research has introduced innovative methods to counteract adversarial attacks.

Frequency Domain Analysis (FDA), a technique that advances the principles of Spectral Signature Matching (SSM). FDA analyzes the frequency components of both input data and gradients, showcasing heightened sensitivity in detecting subtle adversarial perturbations. This method marks a significant improvement over traditional SSM approaches, particularly in identifying less perceptible adversarial attacks[33].

Complementing FDA, Outlier Detection with Autoencoders (ODAE), which employs autoencoders to reconstruct what is considered clean data. Adversarial examples, characterized by significant reconstruction errors, are effectively identified by ODAE. This method emphasizes a data-driven approach in anomaly detection, harnessing the distinct reconstruction capabilities of autoencoders [32].

Another novel approach with Explainable Gradient Consistency (EGC). EGC merges Interpretable Gradient Consistency (IGC) with interpretable saliency maps, thus enabling the identification of specific regions in input data that have been manipulated in adversarial examples. EGC stands out for its transparency and fairness in the detection process, offering visual explanations for identified adversarial inputs. Together, these methods represent significant strides in the ongoing effort to secure machine learning models against sophisticated adversarial threats [22].

Concentrating on adversarial training, the regularization into Graph Neural Networks (GNNs) by considering the inherent structure of the underlying graph. The potential of adversarial training in bolstering the robustness of GNNs. The robustness of GNNs and offers a comprehensive overview

of current research on adversarial attacks, providing valuable insights into both challenges and opportunities for fortifying GNN security. Emphasizing the necessity for collaborative efforts among experts in graph theory, machine learning, and cybersecurity, the study underscores the intricate challenges presented by adversarial attacks on GNNs. Bridging this interdisciplinary gap holds the promise of developing more thorough and effective defense mechanisms [37].

In the evolving landscape of machine learning security, recent studies have introduced innovative approaches to enhance model robustness against adversarial attacks. A method that involves pruning weights that are particularly sensitive to adversarial perturbations during the training phase. This technique aims to improve the model's robustness without incurring a significant loss in accuracy. By selectively eliminating weights that contribute to vulnerabilities, the model becomes more resilient to adversarial manipulations [25].

In a different approach, focused on training models with adversarial examples that are generated from a diverse set of pre-trained models. This strategy significantly enhances the model's ability to generalize and defend against a wide range of unseen adversarial attacks. The diversity in the training process ensures that the model is exposed to a wide spectrum of potential threats, thereby fortifying its defenses [26].

A method that employs an ensemble of models with dynamically adjusted weights. These weights are calibrated based on adversarial confidence scores, which enables the ensemble to adaptively respond to varying degrees of adversarial threats. This method not only improves the robustness of the model but also its adaptability, allowing it to effectively counteract evolving adversarial tactics [10].

Collectively, these studies represents the significant contributions to the field, offering novel strategies to strengthen machine learning models against the continuously advancing nature of adversarial attacks [25][26][31].

The field of adversarial attack mitigation in machine learning continues to evolve with innovative strategies. A method specifically targeting the mitigation of Carlini and Wagner attacks through a technique known as Feature Disentanglement. This approach involves separating the features that are essential for the task prediction from those that are susceptible to adversarial manipulations. By isolating and protecting the vulnerable features, this method effectively counters the sophisticated mechanisms employed in Carlini and Wagner attacks. This separation not only enhances the model's resistance to these specific types of attacks but also maintains the integrity and effectiveness of the model in its primary predictive tasks [22].

In a parallel development, a defense mechanism against DeepFool attacks, employing a technique termed Adaptive Smoothing. This method involves applying a smoothing filter to the input data, which essentially blurs the potential points of attack. By doing so, it becomes significantly more challenging for DeepFool attacks to precisely alter the input data in a way that misleads the model. The key advantage of Adaptive Smoothing is its ability to mitigate attacks without compromising the fidelity of the clean data. This ensures that the model's performance on legitimate data is not adversely affected while enhancing its resilience against these adversarial attacks [22].

Together, the methods developed represents significant advancements in safeguarding machine learning models. They address the dual need of maintaining model accuracy and robustness against increasingly sophisticated adversarial attacks, thus contributing to the overall reliability and security of machine learning systems [8].

The susceptibility of deep learning models lacks emphasis on fostering interdisciplinary collaboration. Closing the gap between machine learning experts, security researchers, and domain-specific professionals is vital for crafting holistic adversarial defense strategies.To address these gaps, the research community needs to delve deeper into the intricate challenges of adversarial attacks. This involves considering diverse application contexts and constructing adaptive, interpretable, and collaborative defense mechanisms. Integration of technical expertise across disciplines is essential for developing comprehensive strategies that mitigate adversarial threats effectively [38].

Utilizing formal verification techniques to mathematically prove the robustness of models against specific attack types offers a promising direction for future research. Incorporating human expertise into detection and mitigation strategies can enhance defense effectiveness, particularly against novel attacks. Evolving Attack Landscape: Continuous adaptation and improvement of defense mechanisms are crucial as attackers develop new and more sophisticated techniques.

A method that focuses on the logit outputs, which are the model's raw predictions before the final activation function like softmax. This method detects adversarial examples by comparing the logit outputs for both the original and the perturbed samples. A significant discrepancy in these logits is indicative of a potential adversarial attack. This approach is particularly effective as it doesn't just rely on the final prediction but probes deeper into the model's processing, making it a more nuanced way to detect subtle adversarial manipulations [1][4].

A defense mechanism leveraging the capabilities of Meshed Tensorflow. This advanced framework is utilized to compute gradients in a way that efficiently detects adversarial examples. The strength of this method lies in its high accuracy, as Meshed Tensorflow allows for a more intricate and detailed analysis of the gradients, which are key to identifying adversarial perturbations [2].

The methods that transform input data into a space where the model exhibits increased robustness to adversarial perturbations. Techniques such as random cropping, color jittering, and various forms of data augmentation are employed to achieve this. These transformations effectively create a more complex training environment, teaching the model to focus on the most relevant features and thereby reducing its sensitivity to adversarial modifications [3].

### B. FGSM Attack

The need for sophisticated defenses against stronger attacks. This includes ensemble methods or robust optimization techniques, which are essential to withstand these advanced adversarial methods.Many defense methods struggle to generalize against novel or varied attack types. It's crucial for

defenses to be evaluated against a diverse array of attacks to ensure their effectiveness in real-world scenarios.Some defense mechanisms require significant computational resources or memory. Balancing efficiency with effectiveness is vital, especially in practical applications where resources may be limited [28].

The field of adversarial defense is characterized as an ongoing race between attackers and defenders. Continuously developing new, robust defense mechanisms is essential to stay ahead of increasingly sophisticated attacks.Choosing the most appropriate defense strategy depends on various factors, including the type of attack, the specific model architecture in use, and the desired balance between accuracy and robustness.While significant strides have been made in developing defenses against adversarial attacks, the field remains dynamic and challenging, with a constant need for innovation and adaptation to new threats.

SSM, focuses on analyzing the power spectrum of both the input and its gradient. Adversarial noise often disrupts the natural data patterns, which can be detected using SSM. This method is particularly effective in identifying subtle noise that deviates from the expected spectral characteristics of legitimate data [6].

Input Gradient Consistency, IGC checks for the consistency of gradients across different input channels. Adversarial manipulations, which typically introduce inconsistencies in these gradients, are effectively flagged by IGC. This method hinges on the premise that legitimate inputs would maintain a certain level of gradient consistency, unlike their adversarial counterparts [7].

Kernel Deep Density Estimators (K-DDEs), learn the underlying distribution of the data and are adept at identifying outliers indicative of adversarial perturbations. This approach is grounded in statistical learning and provides a robust way to detect anomalies that stray from the learned data distribution [8].

A research on DAT involves training the model with a diverse set of adversarial examples. This enhances the model's resilience to future attacks by exposing it to a wide range of potential adversarial tactics during the training phase [9].

### C. Generative Adversarial Attack (GAN)

A method that co-trains the model alongside a GAN, which generates realistic adversarial examples. This joint training enables the model to better distinguish between legitimate and adversarial inputs, thereby improving its robustness [2].

MAT, which utilizes meta-learning algorithms to develop a generalizable strategy for adapting to various types of adversarial attacks. This approach allows models to quickly adapt to new and unseen adversarial tactics based on learned meta-strategies. Each of these methods contributes to a more comprehensive and multi-faceted approach to defending machine learning models against the ever-evolving landscape of adversarial attacks, focusing on both preemptive training and active detection to enhance model robustness and security [11].

A method that utilizes multi-scale gradient filtering to defend against DeepFool attacks. This approach focuses on modifying the gradient information at multiple scales, effectively mitigating the impact of these attacks. A key advantage of this method is its ability to preserve the fidelity of the input data, ensuring that the defensive process does not degrade the quality of legitimate inputs [30].

Certified robustness methods, aiming to guarantee model robustness against adversarial examples within specific norm bounds. These methods provide a mathematical assurance of robustness, offering a more reliable and quantifiable defense against adversarial manipulations [2][10][17][27][31][36].

A novel feature pruning technique to enhance the efficiency of adversarial training. By pruning less relevant features, this technique reduces the computational cost associated with training models on adversarial examples, while still maintaining a high level of robustness against attacks [23][32][37].

The concept of ensemble learning, utilizing a collection of diversified models to improve the detection and mitigation of adversarial examples. This approach bases its defense on the confidence scores from different models, enhancing the overall accuracy and reliability of detecting adversarial attacks [22][24][29][33][38].

Wasserstein distance divergence in the generation of adversarial examples. This method produces more diverse and realistic adversarial inputs for robust training, thereby improving the model's generalizability to unseen attacks. The use of Wasserstein distance helps in creating more challenging and varied adversarial scenarios, which is crucial for comprehensive and effective adversarial training. Each of these studies contributes uniquely to the field of adversarial defense, showcasing the diverse range of approaches being developed to safeguard machine learning models against the continuously advancing techniques of adversarial attacks [34].

Despite significant advancements in adversarial defense mechanisms for machine learning models, there remain challenges in developing universally robust, computationally efficient, and adaptable defense strategies that can effectively counter a wide range of adversarial attacks, including novel and sophisticated ones.

## IV. METHODOLOGY

The MNIST dataset is a collection of grayscale images of handwritten digits (0 through 9). Each image is 28 pixels in height and 28 pixels in width, resulting in a 2D array of pixel values representing the digit.

Let $X \in \mathbb{R}^{M \times N}$ represent an image in the MNIST dataset. Here, $M$ is the height of the image (number of rows), and $N$ is the width of the image (number of columns). Each element $X_{ij}$ of the matrix $X$ corresponds to the intensity value of the pixel at row $i$ and column $j$. The intensity values are real numbers in the range of 0 to 255, where 0 represents black (no intensity) and 255 represents white (maximum intensity). The intensity values are typically integers ranging from 0 to 255, with 0 being completely dark and 255 being fully illuminated. This grayscale representation captures the variations in pixel intensity without considering color information. - $X_{ij}$: Intensity value of the pixel at row $i$ and column $j$. - $M$: Height of the image (number of rows). - $N$: Width of the image (number of columns). - Each image in the MNIST

dataset is essentially a 2D grid of pixels, forming a matrix $X$. - The grayscale intensity values provide information about the darkness or brightness of each pixel. - The size of the matrix ($M \times N$) is fixed for all images in the MNIST dataset (28x28 pixels). - This representation is suitable for machine learning algorithms that can learn patterns and features from the pixel values to recognize handwritten digits.

- If $X_{12} = 150$, it means the pixel at the 1st row and 2nd column has an intensity value of 150, which corresponds to a shade of gray.

Understanding the input data representation is crucial for preprocessing and feeding the data into machine learning models to effectively learn and make predictions based on the patterns within these pixel values.

### A. Feature Masking Process

*1) Setting pixels to constant value:* This can be achieved by setting pixel values in specific regions to a constant value. For example, you can set a rectangular region of the image to 0 or 255. Mathematically:

$$X'_{ij} = \begin{cases} c & \text{if } (i,j) \text{ is in the masked region} \\ X_{ij} & \text{otherwise} \end{cases} \tag{1}$$

Here, $X'_{ij}$ represents the modified pixel value, and $c$ is the constant value.

*2) Applying filters:* Filters, such as blurring or distortion filters, can be applied to certain areas of the image. Let $F$ be a filter matrix, and $*$ denotes the convolution operation. The masked image can be obtained as:

$$X' = X * F \tag{2}$$

*3) Dropout:* Dropout is a technique where certain pixels are randomly set to zero during training. Mathematically:

$$X'_{ij} = X_{ij} \cdot M_{ij} \tag{3}$$

where $M_{ij}$ is a binary mask with elements randomly set to 0 or 1.

### B. Reducing Dependency on Specific Features

Feature masking disrupts the input data in a controlled manner, preventing the neural network from relying too heavily on specific pixels for classification. During feature masking, specific pixels are modified or set to constant values, introducing controlled perturbations to the input data. Mathematically, this disruption is represented by modifying the pixel values, such as in the setting pixels to a constant value or applying filters.The controlled disruption reduces the model's dependence on individual pixel values, promoting a more generalized understanding of the features in the data. By reducing the network's reliance on specific pixels, the model becomes less sensitive to noise or variations in those pixels.This helps the model focus on more relevant features, leading to better generalization on unseen data. Particularly useful in scenarios where certain pixels may be subject to noise or variations that are not indicative of the overall pattern.

### C. Regularization Effect

Feature masking acts as a form of regularization during training, preventing overfitting and encouraging the model to learn more robust and generalizable features. The disruption introduced by feature masking, such as setting pixels to constant values or applying filters, adds noise to the training process. This regularization effect is achieved by modifying the input data in a controlled manner. During dropout, random zeros are introduced in the input, preventing the network from relying too heavily on specific pixel values. Regularization helps prevent overfitting, where the model memorizes training data rather than learning the underlying patterns.Feature masking introduces a level of uncertainty, forcing the model to be more flexible and less prone to memorizing noise.The regularization effect contributes to a more robust model that performs well on unseen data. Feature masking disrupts the input data in a controlled manner, preventing the neural network from relying too heavily on specific pixels for classification. Feature masking acts as a form of regularization during training, preventing overfitting and encouraging the model to learn more robust and generalizable features.

### D. Promoting Invariance

Feature masking encourages the neural network to be invariant to certain changes in the input, making it more resilient to variations in irrelevant features. Let $x$ be the input image represented as a matrix of pixel values. Feature masking is performed by applying a masking function $M(x)$ to $x$.The masking function selectively alters or ignores certain pixel values in $x$, promoting invariance to those specific changes.Mathematically, the result of this operation is represented as:

$$y = M(x)$$

where $y$ is the masked image. Feature masking aims to make the neural network less sensitive to variations in specific regions or features of the input image. The masking function $M(x)$ introduces controlled changes to the input, encouraging the network to focus on more relevant and discriminative features.Invariance to certain changes enhances the model's ability to generalize across different instances of the same class, making it more robust to variations that are irrelevant for classification.

### E. Representation of the Input Image x

Assume the input image $\mathbf{x}$ is represented as a 2D matrix $[x_{ij}]$ where $i$ and $j$ index the rows and columns, respectively. A grayscale image is typically represented as a 2D array of pixel values, where $x_{ij}$ denotes the intensity value of the pixel at row $i$ and column $j$. The grayscale image can be represented as:

$$\mathbf{x} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1N} \\ x_{21} & x_{22} & \dots & x_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M1} & x_{M2} & \dots & x_{MN} \end{bmatrix} \tag{4}$$

Each element $x_{ij}$ represents the intensity value of a pixel in the image. Grayscale images have a single channel, where pixel

values range from 0 (black) to 255(white).This representation is suitable for scenarios where color information is not crucial, such as in the MNIST dataset.

In the case of a color image, $\mathbf{x}$ would typically be a 3D matrix $[x_{ijk}]$, where $k$ indexes the color channel (e.g., RGB channels). A color image is represented as a 3D array, where $x_{ijk}$ represents the intensity value of the pixel at row $i$, column $j$, and color channel $k$. The color image can be represented as:

$$\mathbf{x} = \begin{bmatrix} \begin{bmatrix} x_{111} & x_{112} & \dots & x_{11N} \\ x_{121} & x_{122} & \dots & x_{12N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M11} & x_{M12} & \dots & x_{M1N} \end{bmatrix}, \\ \begin{bmatrix} x_{211} & x_{212} & \dots & x_{21N} \\ x_{221} & x_{222} & \dots & x_{22N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M21} & x_{M22} & \dots & x_{M2N} \end{bmatrix}, \\ \vdots, \\ \begin{bmatrix} x_{1M1} & x_{1M2} & \dots & x_{1MN} \\ x_{2M1} & x_{2M2} & \dots & x_{2MN} \\ \vdots & \vdots & \ddots & \vdots \\ x_{NM1} & x_{NM2} & \dots & x_{NMN} \end{bmatrix} \end{bmatrix} \quad (5)$$

Color images have multiple channels, typically representing Red, Green, and Blue (RGB) color information.The 3D matrix captures color intensity values for each pixel in the image.This representation is essential for tasks where color information plays a crucial role, such as in natural images.

Understanding the representation of input images, whether grayscale or color, involves considering the dimensionality and intensity values associated with each pixel, providing the foundation for image processing and analysis.

*F. The Masking Function $M(\mathbf{x})$*

The masking function $M(\mathbf{x})$ is applied to the image $\mathbf{x}$ and produces a mask matrix of the same dimensions as $\mathbf{x}$. The mask matrix, denoted as $[m_{ij}]$ for grayscale or $[m_{ijk}]$ for color images, is generated by $M(\mathbf{x})$.Each entry in the mask matrix is either 1 or 0, indicating whether to keep or mask the corresponding pixel value in $\mathbf{x}$. For a grayscale image:

$$m_{ij} = \begin{cases} 1 & \text{if } M(\mathbf{x}) \text{ keeps the pixel at } (i, j) \\ 0 & \text{if } M(\mathbf{x}) \text{ masks the pixel at } (i, j) \end{cases} \quad (6)$$

Similarly, for a color image, the mask is represented as $[m_{ijk}]$, where each $m_{ijk}$ is 1 or 0. The masking function is a key element in feature masking processes, such as setting pixels to constant values, applying filters, or using dropout. The binary nature of the mask matrix (1 or 0) signifies the decision to retain or discard pixel information. By controlling which pixels are masked or retained, the masking function influences the model's perception of features during training. Masks can be generated based on different criteria, introducing flexibility in selectively modifying or preserving image elements.

The masking function involves recognizing its role in determining which pixels are retained or masked, providing a mechanism for controlled feature manipulation during various image processing tasks.

*G. Applying the Mask*

The masked image $\mathbf{y}$ is obtained by performing an element-wise multiplication of the original image $\mathbf{x}$ and the mask matrix $M(\mathbf{x})$:

$$\mathbf{y} = \mathbf{x} \odot M(\mathbf{x}) \quad (7)$$

For a grayscale image, the element-wise operation is expressed as:

$$y_{ij} = x_{ij} \times m_{ij} \quad (8)$$

Similarly, for a color image with three channels:

$$y_{ijk} = x_{ijk} \times m_{ijk} \quad (9)$$

In this operation, $m_{ij}$ (or $m_{ijk}$) takes values of 0 or 1. When $m_{ij}$ is 0, the corresponding pixel in $\mathbf{y}$ is effectively masked (set to zero), and when it is 1, the original pixel value is retained. The element-wise masking operation is a fundamental step in feature masking processes, influencing how specific regions or features in the image are modified or retained. When $m_{ij}$ (or $m_{ijk}$) is 0, the corresponding pixel in the resulting masked image $\mathbf{y}$ is suppressed or set to zero. This operation is essential for applying feature-specific modifications, allowing the model to focus on relevant image components while discarding or altering less important ones. The masked image $\mathbf{y}$ retains the structure and features of the original image $\mathbf{x}$ based on the applied masking strategy. The element-wise masking operation provides insights into how feature masking techniques selectively modify or retain pixel values, influencing the learning process of neural networks and other image processing applications.

*H. Purpose and Effects*

Feature masking is a technique commonly employed in image processing and deep learning to direct a model's attention to specific regions of an image, augment data, or simulate occlusions during training for increased robustness.

Feature masking simplifies to selectively zeroing out certain pixels while leaving others unchanged. This is represented by the element-wise multiplication of the original image $\mathbf{x}$ and the mask matrix $M(\mathbf{x})$:

$$\mathbf{y} = \mathbf{x} \odot M(\mathbf{x}) \quad (10)$$

For grayscale images:

$$y_{ij} = x_{ij} \times m_{ij} \quad (11)$$

For color images:

$$y_{ijk} = x_{ijk} \times m_{ijk} \quad (12)$$

Feature masking alters the input data fed into the model by selectively modifying pixel values based on the mask. When certain pixels are zeroed out (masked), the model focuses on the remaining unmasked pixels during training. This allows the model to learn features that are relevant for classification or other tasks while ignoring or being less sensitive to specific regions.

Feature masking directs the model's focus to specific features or regions, enabling it to learn discriminative patterns.

Useful in scenarios where certain image components are more critical for decision-making. By selectively altering pixels, feature masking contributes to data augmentation, introducing variations in the training data. This helps improve the model's generalization by exposing it to diverse instances of the same class. Simulating occlusions during training with feature masking enhances the model's robustness to partial or obscured input images.The model learns to make predictions even when parts of the input are hidden or occluded.

The purpose and mathematical significance of feature masking provides a powerful tool for shaping the learning process of models, enhancing their ability to generalize, and improving robustness to variations in input data.

### I. Feature Masking as Dimensionality Reduction

Consider a neural network with input features represented by a vector $X \in \mathbb{R}^d$, where $d$ is the original dimensionality of the input space. The feature masking process involves element-wise multiplication of the input features by a binary mask $M \in \{0,1\}^d$ that determines which features are active (1) or masked (0). The masked input $\tilde{X}$ can be mathematically expressed as:

$$\tilde{X} = M \odot X \qquad (13)$$

$X \in \mathbb{R}^d$ represents the original input features, where $d$ is the dimensionality of the input space. $M \in \{0,1\}^d$ is a binary mask vector, indicating which features are active (1) and which are masked (0). The element-wise multiplication $\odot$ is performed between the input features $X$ and the binary mask $M$, resulting in a masked input $\tilde{X}$. The element-wise multiplication is expressed as:

$$\tilde{X}_i = M_i \times X_i \qquad (14)$$

where $i$ represents the index of each feature in the vectors.

The binary mask $M$ provides control over which features are allowed to contribute to the neural network's computations. Active features (where $M_i = 1$) retain their original values, while masked features (where $M_i = 0$) are effectively set to zero. Feature masking can be seen as a form of dimensionality reduction, as it allows the network to focus on a subset of relevant features. This is particularly useful when certain features are noisy or irrelevant to the learning task. Feature masking acts as a form of regularization by introducing sparsity in the input space. Sparse inputs encourage the neural network to learn more robust and generalizable features.

The feature masking process in a neural network involves selectively modifying input features based on a binary mask, influencing the model's attention, reducing dimensionality, and providing regularization. This process is valuable for enhancing the network's ability to learn meaningful patterns from the input data.

### J. Regularization Objective

Regularization is often expressed through an additional term in the loss function. In the case of feature masking, the regularization term encourages sparsity in the mask, penalizing the model for relying too much on specific features. The overall loss function $(L)$ can be written as a combination of the standard task-specific loss $(L_{\text{task}})$ and a regularization term $(R)$:

$$L = L_{\text{task}} + \lambda R \qquad (15)$$

$L_{\text{task}}$ represents the standard task-specific loss, measuring the model's performance on the primary learning task. $R$ is the regularization term, which penalizes the model for non-ideal behaviors, such as relying too heavily on specific features. $\lambda$ is a hyperparameter that controls the strength of the regularization. It determines how much importance is given to the regularization term relative to the task-specific loss. The combination of $L_{\text{task}}$ and $\lambda R$ creates a trade-off: the model aims to minimize the task-specific loss while keeping the regularization term in check. The regularization term $(R)$ associated with feature masking might involve measuring the sparsity of the mask:

$$R = \sum_{i=1}^{d} |M_i| \qquad (16)$$

where $d$ is the dimensionality of the input features. The regularization term encourages sparsity in the mask by penalizing non-zero entries. This leads to feature selection, allowing the model to focus on a subset of relevant features. The hyperparameter $\lambda$ controls the trade-off between minimizing task-specific loss and minimizing the impact of the regularization term. A higher $\lambda$ encourages stronger regularization, limiting the model's reliance on specific features. By penalizing the model for overfitting to certain features, feature masking regularization improves the generalization capability of the model. The model becomes less sensitive to noise or irrelevant features in the input.

The incorporation of feature masking regularization in the loss function provides a mechanism for controlling the sparsity of the mask, balancing task-specific learning with the encouragement of more generalized feature dependencies. This regularization contributes to building models that generalize well to new and unseen data.

### K. Training with Masks

During training, different masks are applied to the input data in a stochastic manner. This can be represented as a probability distribution over masks. Let $P(M)$ be the probability distribution of masks, and $E_M$ denote the expectation over masks. The training objective can be expressed as the minimization of the expected loss:

$$\min_{\theta} E_M[L(f(X \odot M; \theta), y)] + \lambda R(M) \qquad (17)$$

$P(M)$ represents the probability distribution over masks. Each mask $M$ is a realization from this distribution during training. - $E_M$ is the expectation operator over masks, indicating that the training objective involves averaging the loss over different mask realizations. $L(f(X \odot M; \theta), y)$ is the task-specific loss, measuring the model's performance on the primary learning task with a masked input. $R(M)$ is the regularization term that penalizes non-ideal behaviors, such as sparsity in the mask. $\lambda$ is a hyperparameter controlling the trade-off between the task-specific loss and the regularization term. The overall objective is to minimize the

expected loss by considering the variability introduced by different masks during training. Stochastic masking introduces randomness during training by applying different masks in a probabilistic manner. This randomness helps the model generalize better and become more robust to variations in input data. The expectation $E_M$ represents the average loss over all possible masks, capturing the model's performance under diverse feature conditions. This averaging helps mitigate the impact of individual masks that may be overly specific or noisy. The regularization term $R(M)$ encourages certain properties in the mask distribution, such as sparsity. This helps prevent the model from overfitting to specific features and promotes a more generalized understanding of the input. The stochastic feature masking during training involves considering the variability introduced by different masks, the expectation over masks, and the joint optimization of task-specific loss and regularization. This approach contributes to the model's ability to adapt to diverse input conditions and enhances its overall robustness.

*L. Adversarial Robustness*

The concept of adversarial robustness can be framed in terms of the impact of perturbations on the masked input space. If $X_{\text{adv}}$ is an adversarial perturbation added to $X$, the masked adversarial input $\tilde{X}_{\text{adv}}$ can be expressed as:

$$\tilde{X}_{\text{adv}} = M \odot (X + X_{\text{adv}}) \tag{18}$$

$X_{\text{adv}}$ represents the adversarial perturbation added to the original input $X$. The masked input space is modified by element-wise multiplication $\odot$ with the binary mask $M$. The expression $X + X_{\text{adv}}$ represents the addition of the original input and the adversarial perturbation. The mask $M$ selectively applies perturbations to certain features, influencing the impact of adversarial perturbations. The resulting $\tilde{X}_{\text{adv}}$ is the masked adversarial input.

The element-wise multiplication with the mask $M$ allows for selective application of perturbations to the input features. Certain features, determined by the mask, may be more or less susceptible to adversarial perturbations. The mask $M$ plays a crucial role in shaping the adversarial robustness of the model. By controlling which features are affected by perturbations, the mask contributes to the model's resilience against adversarial attacks. Understanding the impact of adversarial perturbations in the masked input space helps in developing models that generalize well in the presence of adversarial examples. The model learns to be robust to variations introduced by adversarial perturbations while focusing on relevant features. Framing adversarial robustness in the context of the masked input space involves considering how perturbations selectively impact features based on the binary mask, influencing the model's resilience against adversarial attacks. This approach contributes to the development of more robust machine learning models.

*M. Mathematical Framework of Feature Masking and Data Augmentation*

- Description: The stochastic application of masks during training is a form of data augmentation. Mathematically, data augmentation introduces variability in the

training data to improve generalization. In the context of feature masking, variability is directly injected into the feature space through different masks, encouraging the model to generalize better to diverse input patterns.

- Mathematical Significance:
  - Data Augmentation as Variability Introduction: Data augmentation is represented mathematically by introducing variability in the training data. In feature masking, this variability is introduced directly into the feature space through the application of different masks during training. Mathematically, data augmentation can be seen as modifying the input data $X$ through a stochastic process:

    $$X' = \text{Augmentation}(X)$$

  - Feature Masking Mathematical Framework: Feature masking involves masks, regularization terms, and expectations over mask distributions. During training, the masked input $\tilde{X}$ is obtained by element-wise multiplication with a mask:

    $$\tilde{X} = M \odot X$$

    The regularization term $R(M)$ encourages sparsity in the mask to prevent overreliance on specific features. Expectations over mask distributions are incorporated into the training objective:

    $$\min_{\theta} E_M[L(f(\tilde{X}; \theta), y)] + \lambda R(M)$$

- **Insights:**
  - Diversity in Features for Decision-Making: Feature masking encourages diversity in the features used for decision-making during training. By applying different masks stochastically, the model learns to be invariant to variations in the input. This diversity enhances generalization by exposing the model to a broader range of input patterns.
  - Formalization of Regularization: The regularization term $R(M)$ ensures that the model does not overly rely on specific features, promoting more robust and generalized learning. The regularization effect is formalized in the loss function, contributing to improved model performance on unseen data.
  - Alignment with Adversarial Robustness: Feature masking, by controlling the impact of perturbations through masks, aligns with principles of adversarial robustness. The model learns to be resilient to adversarial attacks by considering diverse feature spaces.

*1) Random masking as a stochastic process:* Consider the training images as a set $\{x^{(1)}, x^{(2)}, \ldots, x^{(n)}\}$, where each $x^{(i)}$ is an image. A random mask $M^{(i)}$ is applied to each image during each epoch of training, which can be mathematically represented as a stochastic process. The masked image is then $M^{(i)}(x^{(i)})$, where the operation $M^{(i)}$ selectively alters pixel

values in $x^{(i)}$ based on a random pattern. Let $X$ represent the set of training images: $X = \{x^{(1)}, x^{(2)}, \ldots, x^{(n)}\}$. A random mask $M^{(i)}$ is applied to each image $x^{(i)}$ during training epochs. This can be expressed as a stochastic process:

$$M^{(i)}(x^{(i)}) = M^{(i)} \odot x^{(i)} \qquad (19)$$

Here, $M^{(i)}$ is a binary mask, and $\odot$ denotes element-wise multiplication. The stochastic process introduces variability in the training data by applying different masks to each image during each epoch. Random masking as a stochastic process introduces variability in the training data. This variability arises from the different random masks applied to each image during each training epoch. The element-wise multiplication ($\odot$) selectively alters pixel values in $x^{(i)}$ based on the binary mask $M^{(i)}$. The process results in a diverse set of masked images for each input in the training set. This diversity promotes a richer learning experience for the model by exposing it to various instances of the same image with different masked patterns. By training on a dataset with masked images generated through a stochastic process, the model becomes more robust to variations in input patterns. The introduction of variability enhances the model's ability to generalize and make predictions on unseen data.

*2) Training with masked inputs:* In the training of neural networks, rather than learning a mapping $f(x^{(i)})$ directly, a stochastic masking process is incorporated. Each training image $x^{(i)}$ undergoes modification through a random mask $M^{(i)}$, resulting in $M^{(i)}(x^{(i)})$. The neural network learns a mapping $f(M^{(i)}(x^{(i)}))$ during training. Here, $x^{(i)}$ represents the matrix of pixel values, and $M^{(i)}(x^{(i)})$ is another matrix with modified entries based on the applied mask.

- Let $X$ denote the set of training images: $X = \{x^{(1)}, x^{(2)}, \ldots, x^{(n)}\}$.

- The stochastic masking process is represented mathematically as:
$$f(M^{(i)}(x^{(i)}))$$

- The application of $M^{(i)}$ to $x^{(i)}$ modifies each entry of the matrix element-wise, enforcing a focus on different features in each iteration:
$$M^{(i)}(x^{(i)}) = M^{(i)} \odot x^{(i)}$$

- The neural network adapts to variations introduced by the stochastic masking process, resulting in a mapping that is inherently robust and less prone to overfitting.

- Feature Variation in Training:
  - The alteration induced by $M^{(i)}$ forces the neural network to focus on different features of the input in each iteration.
  - This variation in training instances helps prevent the network from over-relying on specific features, contributing to improved generalization.

- Enhanced Robustness:
  - The network's exposure to $M^{(i)}(x^{(i)})$ during training promotes adaptability to variations in input patterns.

  - This enhanced robustness makes the network more capable of handling diverse inputs, leading to improved performance on unseen data.

- Prevention of Overfitting:
  - Stochastic masking serves as a regularization technique by introducing variability in the training process.
  - This variability prevents the network from memorizing specific details in the training data, reducing the risk of overfitting to noise.

- Improved Generalization:
  - By learning a mapping $f(M^{(i)}(x^{(i)}))$ instead of $f(x^{(i)})$, the network becomes more adept at generalizing its knowledge to novel instances.
  - The focus on diverse features through stochastic masking contributes to a model that can better handle different variations in the input space.

*3) Consistency in testing:* During the testing phase, the input $x_{\text{test}}$ is subjected to two scenarios: either it is not masked at all, or a consistent mask $M_{\text{test}}$ is applied. The model's performance is evaluated using $f(M_{\text{test}}(x_{\text{test}}))$ or $f(x_{\text{test}})$, ensuring a consistent and fair evaluation. This approach maintains control over testing conditions, allowing for a clear comparison of the model's performance with and without masking.

- During testing, the evaluation is carried out under two conditions:
  - Without masking: $f(x_{\text{test}})$
  - With consistent masking: $f(M_{\text{test}}(x_{\text{test}}))$

- The application of $M_{\text{test}}$ to $x_{\text{test}}$ follows a similar mathematical representation as in the training phase:
$$M_{\text{test}}(x_{\text{test}}) = M_{\text{test}} \odot x_{\text{test}}$$

- This consistency ensures that the model is tested under controlled conditions, allowing for a fair and unbiased assessment of its performance.

- Controlled Evaluation:
  - By evaluating the model under two distinct conditions (with and without masking), consistency in testing provides a controlled environment for performance assessment.
  - This controlled evaluation is crucial for understanding how well the model generalizes to both unaltered and consistently masked inputs.

- Fair Model Comparison:
  - Consistent testing enables a fair comparison of the model's performance under different conditions.
  - This comparison is valuable in assessing the impact of stochastic masking on the model's predictions and understanding its robustness to variations introduced during training.

- Understanding Masking Influence:
  - Testing with and without masking allows for a clear understanding of how the stochastic

masking process influences the model's behavior during inference.

- Insights gained from consistent testing contribute to refining the model and optimizing its performance for diverse scenarios.

- Robustness Validation:
  - Evaluating the model on both masked and unmasked inputs serves as a validation of its robustness.
  - The consistent testing approach ensures that the model's performance is not skewed by the presence or absence of masking during evaluation.

*4) Hyperparameter tuning:* The design of the mask $M$ is a critical hyperparameter, involving the proportion of features masked, the pattern of masking, and the variability between epochs. Mathematically, this can be viewed as tuning the parameters of the stochastic process governing $M$, balancing the network's exposure to features with the need for robustness and generalization.

- Stochastic Process $M$:
  - Let $M$ represent the stochastic process of masking during training.
  - The application of $M$ to an input $x^{(i)}$ is given by $M(x^{(i)}) = M \odot x^{(i)}$, where $\odot$ denotes element-wise multiplication.

- Hyperparameters of $M$:
  - Designing $M$ involves tuning hyperparameters that govern the stochastic process:
    - Proportion of features masked.
    - Pattern of masking.
    - Variability between epochs.

- Mathematical Tuning:
  - The design process can be expressed mathematically as the tuning of hyperparameters:

$$M = \text{Tune}(p, \text{pattern}, \text{variability}) \quad (20)$$

  where $p$ is the proportion of features masked, pattern specifies the masking pattern, and variability controls the variability between epochs.

- Trade-off between Diversity and Consistency:
  - The hyperparameters influence the trade-off between diversity and consistency in the training process.
  - A higher $p$ introduces more diversity by masking a larger proportion of features, while a lower $p$ maintains consistency.
  - The masking pattern and variability further contribute to this balance.

- Exposure to Features:
  - Adjusting hyperparameters allows control over the network's exposure to features. Higher values of $p$ promote increased variability, exposing the model to a broader range of input patterns.

- Robustness and Generalization:
  - Tuning the hyperparameters impacts the model's robustness and generalization capabilities. Striking the right balance ensures that the model can adapt to diverse inputs while maintaining consistency.

- Trade-off Considerations:
  - The proportion of features masked ($p$) serves as a key trade-off parameter. A delicate balance is needed to prevent overfitting (too much diversity) or underfitting (too much consistency).

- Pattern and Variability Impact:
  - The choice of masking pattern and variability between epochs contributes to the richness of the training data. Patterns that capture relevant features and controlled variability enhance the learning process.

- Iterative Tuning:
  - The design of $M$ involves an iterative tuning process. Hyperparameters may be adjusted based on the network's performance, ensuring a dynamic adaptation to the learning dynamics.

*5) Regularization and reduced dimensionality:* From a regularization standpoint, feature masking can be seen as adding a form of noise to the input data, which helps in preventing overfitting. Mathematically, this reduces the effective dimensionality of the input space, as the network is forced to make predictions with incomplete information, enhancing its ability to generalize.

- Feature Masking Operation:
  - Let $x^{(i)}$ represent the input data. The feature masking operation is defined as:

$$x_{\text{masked}}^{(i)} = M^{(i)} \odot x^{(i)} \quad (21)$$

  - Here, $M^{(i)}$ is a binary mask, and $\odot$ denotes element-wise multiplication.

- Regularization Effect:
  - The feature masking introduces noise by selectively setting certain features to zero, creating an incomplete representation of the input during training.
  - Mathematically, this can be expressed as injecting randomness into the input data:

$$x_{\text{masked}}^{(i)} = \text{RandomMask}(x^{(i)}) \quad (22)$$

- Reduced Effective Dimensionality:
  - The masking operation reduces the effective dimensionality of the input space. It limits the information available to the network for each instance during training.
  - Mathematically, this reduction can be quantified as:

$$\text{Effective Dimensionality} = \sum_{j=1}^{D} m_j \quad (23)$$

where $D$ is the original dimensionality, and $m_j$ is the binary value of the $j$-th element in the mask.

- Noise Introduction for Regularization:
  - Feature masking introduces noise by hiding certain features during each training instance.
  - This noise prevents the model from memorizing specific patterns, promoting a more generalized understanding of the data.

- Preventing Overfitting:
  - The regularization effect of feature masking helps in preventing overfitting by discouraging the model from relying too heavily on specific details present in the training data.
  - The network learns to make predictions with a more robust understanding of the underlying patterns.

- Generalization Enhancement:
  - By training on partially masked data, the network becomes more adept at generalizing to unseen instances.
  - The reduced effective dimensionality forces the model to focus on essential features, improving its ability to generalize to diverse inputs.

- Adaptability to Incomplete Information:
  - Feature masking encourages the model to be adaptable to incomplete information, mimicking real-world scenarios where not all features may be available during prediction.
  - This adaptability contributes to the model's resilience and performance on diverse datasets.

*6) Robustness against adversarial attacks:* Adversarial attacks often exploit specific weaknesses in the model's learned mapping $f(x)$. By training the network on $f(M(x))$, where $M$ varies, the model becomes less sensitive to specific patterns and more resilient to such manipulations.

$$f(M(x))$$

The variability introduced by the stochastic masking process reduces the model's reliance on specific features, making it more robust against adversarial attacks targeting those features.

- Adversarial Mapping:
  - Let $f(x)$ represent the learned mapping of the network on clean data.
  - Adversarial attacks often aim to exploit vulnerabilities in $f(x)$ by manipulating input patterns.

- Stochastic Masking Operation:
  - The network is trained on $f(M(x))$, where $M$ is a stochastic mask applied to the input data $x$.
  - Mathematically, this can be expressed as:

$$f(M(x))$$

- Variability in Training:
  - The stochastic masking process introduces variability in the training data by applying different masks to each input during each training instance.
  - The variability is controlled by the stochastic mask $M$, leading to diverse instances of masked inputs.

- Robustness against Adversarial Attacks:
  - The introduced variability reduces the model's sensitivity to specific patterns in the input, making it less susceptible to adversarial attacks targeting those patterns.
  - Adversarial attacks crafted for specific features are less effective when the model is trained on $f(M(x))$ due to the unpredictable variations introduced by different masks.

- Reduced Sensitivity to Specific Patterns:
  - Training on $f(M(x))$ introduces unpredictability in the training data, reducing the model's reliance on specific features during inference.
  - This reduced sensitivity makes the model more robust against adversarial attacks that target specific patterns in the input.

- Enhanced Generalization to Varied Inputs:
  - The variability introduced by stochastic masking enables the model to generalize better to a diverse set of inputs.
  - This enhanced generalization contributes to the model's ability to handle variations introduced by adversarial attacks.

- Resilience to Manipulations:
  - Adversarial attacks typically manipulate inputs in a way that exploits the model's vulnerabilities.
  - Training on $f(M(x))$ makes the model more resilient by diminishing the effectiveness of attacks focused on specific patterns.

- Dynamic Defense Mechanism:
  - Stochastic masking serves as a dynamic defense mechanism, making it challenging for adversaries to craft universal attacks that consistently succeed across different instances of the same input.

## V. ALGORITHM

1) Initialization
   Input: Training dataset (e.g., MNIST dataset), neural network model.
   Parameters: Masking ratio $r$ (proportion of features to mask), masking pattern (random or fixed), number of epochs $E$, learning rate $\eta$, batch size $B$.
2) Preprocessing
   Normalize the dataset: Scale the pixel values to a range (e.g., 0 to 1).
   Split the dataset into training and validation sets.
3) Mask Generation
   Define a function `generate_mask(image_shape,`

`ratio`) that creates a mask for an image. The mask should have the same dimensions as the image.

If using random masking, this function generates a new mask for each image in each epoch.

For fixed masking, generate a predefined mask and apply it consistently.

4) Training Loop

For each epoch $e$ in $\{1, 2, \ldots, E\}$:

    a) Shuffle the training dataset.

    b) For each mini-batch $b$ in the training dataset:

       i) For each image $x_i$ in the mini-batch:

          A) Generate a mask $M_i$ using `generate_mask`.

          B) Apply the mask: $\tilde{x}_i = x_i \odot M_i$, where $\odot$ denotes element-wise multiplication.

          C) Perform a forward pass with the masked inputs $\tilde{x}_i$.

          D) Compute the loss $L$ (e.g., cross-entropy loss for classification).

          E) Backpropagate the error and update the model parameters using an optimizer (e.g., SGD, Adam) with a learning rate $\eta$.

5) Validation: After each epoch, evaluate the model on the validation set without applying feature masking. Monitor performance metrics like accuracy, loss, etc.

6) Hyperparameter tuning: Optionally, perform hyperparameter tuning for $r$, $\eta$, and $B$ based on validation performance.

7) Model evaluation: After training, evaluate the final model on a separate test set.
Compare the performance with and without feature masking to assess the impact.

8) Deployment: Deploy the trained model for inference. Optionally, use consistent feature masking if it was part of the training.

## VI. EXPERIMENTAL SET UP

This research investigates the effectiveness of feature masking as a defensive technique against adversarial attacks on neural networks, specifically focusing on the MNIST dataset. The study comprises several key phases, each contributing to a comprehensive evaluation of the proposed approach. We established a baseline by training a standard neural network architecture on MNIST without feature masking, followed by implementing a feature masking algorithm and systematically testing its impact on model performance. Adversarial attacks were simulated using popular methods like the Fast Gradient Sign Method (FGSM) and Projected Gradient Descent (PGD). The experiment configuration and parameters are detailed below to ensure repeatability.

*1) Data preparation:* Dataset: MNIST dataset comprising 60,000 training images and 10,000 test images. Image Characteristics: 28x28 pixel grayscale images of handwritten digits (0 to 9). Pixel Normalization: Scale pixel values from 0 to 255 to a range of 0 to 1. Dataset Splitting: Training set for model training, validation set for hyperparameter tuning, and test set for unbiased model evaluation.

*2) Model architecture:* Neural Network Types: Simple Convolutional Neural Network (CNN) and Multi-Layer Perceptron (MLP). Convolutional Layers: One or two layers with ReLU activation. Pooling Layers: Follow each convolutional layer with max pooling. Fully Connected Layers: One or two layers for classification, with 10 neurons in the final layer using softmax activation. Hidden Layers: One or more hidden layers (e.g., 128 or 256 neurons) with ReLU activation. Flattening: Flatten 28x28 images into a 784-dimensional vector for input. Consistency: Maintain consistent architecture across models for fair comparison.

*3) Training configuration:* Hyperparameters: Keep learning rate, batch size, and number of epochs consistent. Regularization: Depending on model performance, consider dropout or L2 regularization to prevent overfitting.

*4) Feature masking experiment:* Baseline Model: Train a neural network without feature masking, record accuracy, and loss metrics on the test set. Feature Masking Algorithm: Implement a feature masking algorithm and apply various masks to training images across epochs. Consistent Architecture: Ensure the masked model maintains the same architecture as the baseline for fair comparison. Masking Ratios and Patterns: Experiment with different masking ratios and patterns (random to fixed) to determine optimal masking strategy.

*5) Adversarial attack simulation:* Adversarial Methods: Use FGSM and PGD to simulate adversarial attacks on the MNIST test set. Testing: Evaluate both baseline and feature-masked models for robustness against adversarial manipulation.

*6) Result analysis:* Performance Metrics: Assess accuracy and loss on the test set for baseline and feature-masked models. Adversarial Robustness: Analyze model performance under simulated adversarial attacks. By documenting the detailed experimental setup and parameters, we aim to provide a foundation for reproducibility and further exploration of feature masking as a viable strategy for enhancing adversarial defense in neural networks.

To experimentally evaluate the effectiveness of feature masking in enhancing adversarial defense for neural networks, specifically on the MNIST dataset, you need to set up a controlled experiment. This setup will involve comparing the performance of a neural network trained with feature masking against one trained without it, under various conditions.

When working with the MNIST dataset in a machine learning context, the process typically involves two main stages: data preparation and defining the model architecture. Here's a detailed description of each stage:

MNIST dataset is a classic in the field of machine learning, particularly for image recognition tasks. It contains 60,000 training images and 10,000 test images. Image Characteristics: Each image in the MNIST dataset is a 28x28 pixel grayscale image of a handwritten digit (ranging from 0 to 9).

The pixel values in each image, which originally range from 0 to 255, should be normalized to a range of 0 to 1. This is done by dividing each pixel value by 255. Normalization helps in speeding up the training process by ensuring that all input features (pixel values) are on a similar scale. Dataset Splitting:

The dataset should be divided into three subsets: training, validation, and test sets. The training set is used for training the model. The validation set is used to tune hyperparameters and to provide an unbiased evaluation of a model fit during the training phase. The test set is used to provide an unbiased evaluation of the final model fit. For the MNIST dataset, both a simple Convolutional Neural Network (CNN) and a Multi-Layer Perceptron (MLP) can be effective. The choice depends on the complexity of the model you wish to use and the computational resources available. Convolutional Layers: Begin with one or two convolutional layers. These layers extract features from the images by sliding a filter across the input. Each convolutional layer is typically followed by a non-linear activation function like ReLU (Rectified Linear Unit). Pooling Layers: Follow each convolutional layer with a pooling layer (like max pooling) to reduce the spatial size of the representation, reducing the number of parameters and computation in the network. Fully Connected Layers: After the convolutional and pooling layers, add one or two fully connected layers for classification. The last fully connected layer should have 10 neurons (corresponding to the 10 digits) and use a softmax activation function to output probabilities for each digit. Flatten the 28x28 images into a 784-dimensional vector to serve as the input layer. Hidden Layers: Have one or more hidden layers with a sufficient number of neurons (e.g., 128 or 256). Use ReLU for the activation function. The final layer should be a fully connected layer with 10 neurons (one for each digit) with a softmax activation function for classification. To ensure fair comparison in experiments, it's crucial to keep the model architecture consistent. This means using the same number of layers, the same number of neurons in each layer, and the same activation functions. Hyperparameters like learning rate, batch size, and number of epochs should also be kept consistent, unless the specific experiment involves varying these parameters. Depending on the model's performance, regularization techniques like dropout or L2 regularization can be used to prevent overfitting.

In the study we systematically investigated the effectiveness of feature masking as a defensive technique against adversarial attacks on the MNIST dataset. The research was structured into several key phases, each contributing to a comprehensive evaluation of the proposed approach.

Initially, we established a baseline by training a standard neural network architecture on the MNIST dataset without feature masking. This baseline model's performance metrics, notably accuracy and loss on the test set, provided a reference point for subsequent comparisons. Following this, we implemented a feature masking algorithm, applying various masks to the training images across epochs. The neural network, consistent in architecture with the baseline model, was then trained on this modified dataset. This phase included experimentation with different masking ratios and patterns, ranging from random to fixed masking, to ascertain the optimal masking strategy.

Further, we simulated adversarial attacks using prevalent methods such as the Fast Gradient Sign Method (FGSM) and Projected Gradient Descent (PGD), generating adversarial examples from the MNIST test set. These examples were used to test both the baseline and feature-masked models, allowing us to assess their respective robustness to adversarial manipulation.

The performance of both models was meticulously compared using standard metrics like accuracy, precision, recall, and F1-score, on both the normal and adversarial test sets. This comparison provided crucial insights into the effectiveness of feature masking in enhancing model robustness. Additionally, hyperparameter tuning, focusing on aspects such as masking ratio, learning rate, and the number of training epochs, was conducted, utilizing the validation set performance for guiding tuning decisions.

Finally, we conducted statistical tests, to ascertain the significance of the differences observed in the performance metrics between the baseline and feature-masked models. This statistical analysis was pivotal in ensuring the reliability and validity of our findings.

Our research contributes to the growing body of knowledge in neural network security, providing evidence that feature masking can be an effective strategy in augmenting the robustness of neural networks against adversarial attacks. This approach, particularly suitable for simple input domains like MNIST, signifies a strategic advancement in defensive machine learning methodologies. The core of our analysis involved comparing the performance metrics - accuracy, loss, precision, recall, and F1-score - of models trained with and without feature masking. This comparative study was crucial in highlighting the differences in model performance on both standard and adversarially perturbed test sets, thereby providing a clear measure of the effectiveness of feature masking in standard and adversarial contexts. A significant aspect of our research focused on analyzing how varying masking ratios and patterns, such as random versus fixed masking, influenced the model's overall robustness and performance. This analysis was instrumental in identifying the optimal masking strategy, providing valuable insights into the balance between model exposure to features and its ability to generalize and withstand adversarial manipulation. In our pursuit of a more nuanced understanding, we conducted additional experiments to test the model's resilience against a variety of adversarial attack types and strengths. This helped in ascertaining the breadth of the model's robustness. Furthermore, we experimented with combining feature masking with other defense techniques, assessing whether such integrations could further enhance model robustness.

The research necessitated substantial computational resources, with an emphasis on the use of GPUs for expedited training and evaluation. We employed advanced machine learning frameworks like TensorFlow and PyTorch for model development, alongside libraries such as CleverHans and Foolbox for generating a diverse array of adversarial examples. Additionally, we were mindful of the accuracy-robustness trade-off, often observed in adversarial defense mechanisms. We also documented resource utilization to provide insights into the practical feasibility of our methods. The ethical implications of our research were also a paramount consideration, particularly in terms of the potential for adversarial knowledge misuse. We emphasized responsible use and communication of our findings, underlining the importance of advancing AI security in a conscientious manner. Our research provides substantial evidence supporting the effectiveness of feature masking in bolstering neural network security against adversarial threats.

Our comprehensive and structured experimental setup offers valuable insights into the strengths and limitations of feature masking, contributing significantly to the field of neural network security and robust machine learning.

## VII. Experimental Results

The primary objective of the study was to enhance the robustness of a baseline model, which consisted of either a Multilayer Perceptron (MLP) or a Simple Convolutional Neural Network (CNN). Although the baseline model exhibited high accuracy (98%) on the test set, it was found to be susceptible to adversarial attacks. The study aimed to investigate the impact of incorporating feature masking into the model architecture as a means to improve its robustness.

The baseline model was constructed using either a Multilayer Perceptron (MLP) or a Simple Convolutional Neural Network (CNN). Both configurations achieved a baseline accuracy of 98% on the test set. Despite this high accuracy, the baseline model displayed vulnerability to adversarial attacks, leading to performance degradation in the presence of perturbations.

The feature-masked model retained the same architecture as the baseline model. The key modification involved the incorporation of feature masking during training. Feature masking is a technique where a certain percentage of input features is randomly masked or set to zero during each training epoch. The study experimented with different masking ratios, specifically 10%, 30%, and 50%, and applied random masking in each epoch.

During training, the feature-masked model underwent multiple epochs, and in each epoch, a portion of input features was randomly masked based on the specified ratio. This dynamic masking approach aimed to enhance the model's adaptability and robustness by preventing it from relying too heavily on specific features.

The study observed that the accuracy of the feature-masked model on the test set varied with the masking ratio. Specifically, the accuracy decreased from 96% with 10% masking to 92% with 50% masking. This reduction in accuracy can be attributed to the loss of information due to feature masking. However, the primary focus was on the model's robustness to adversarial attacks.

In contrast to the baseline model, the feature-masked model exhibited significantly higher robustness to adversarial attacks. Adversarial attacks typically involve introducing perturbations to the input data to mislead the model. The feature-masked model, despite the reduction in overall accuracy with increased masking ratios, showed less performance degradation under adversarial conditions. This indicates that the model was able to maintain a higher level of performance in the presence of perturbations, showcasing the effectiveness of the feature masking technique in enhancing robustness.

This study demonstrated that the incorporation of feature masking in a neural network model, despite a marginal decrease in accuracy on clean data, can lead to a substantial improvement in robustness to adversarial attacks. This finding has implications for deploying models in real-world scenarios where resilience to adversarial inputs is crucial for reliable performance.

In Fig. 1, the training loss, training accuracy, and validation accuracy are depicted. The figure illustrates the evolution of these metrics throughout the training process.



Fig. 1. Training loss, training and validation accuracy.

In Fig. 2, the comparison between the original and adversarial images is presented.

In Fig. 3, the model's performance is illustrated in the context of both original and adversarial images.

In Fig. 4, the precision, recall, and F1 score trends over epochs are depicted.

In Fig. 5, the confusion matrix provides a visual representation of the model's classification performance.

In Table I, the evaluation results depict the impact of feature masking on model robustness under FGSM attacks.

In Table II, the reported values represent various performance metrics of the model.

In our comprehensive analysis of masking parameters, a crucial trade-off was identified. Specifically, as the masking ratio increased, the model's robustness against adversarial attacks improved, but at the expense of a reduction in overall accuracy. For example, employing a 10% masking ratio resulted in a minor accuracy decrease compared to the baseline, yet it significantly enhanced the model's resistance to adversarial attacks. Conversely, a 50% masking ratio yielded the highest level of robustness but at the cost of a more pronounced accuracy loss. This observation emphasizes the imperative of striking a balance between accuracy and security, tailoring the choice of masking ratio to the specific requirements of the application in question.

TABLE I. Evaluation of Model Robustness with Feature Masking for FGSM Attack

| Model Type | Masking Ratio | Masking Pattern | Training Accuracy | Test Accuracy | Robustness (Accuracy Under Attack) | Training Time Increase |
|---|---|---|---|---|---|---|
| Baseline (No Masking) | N/A | N/A | 99% | 98% | 60% | 0% |
| Feature Masked | 10% | Random | 98% | 96% | 75% | 5% |
| Feature Masked | 30% | Random | 97% | 94% | 80% | 10% |
| Feature Masked | 50% | Random | 95% | 92% | 85% | 15% |
| Feature Masked (Fixed) | 30% | Fixed | 97% | 93% | 78% | 7% |



Fig. 2. Original and adversarial image.



Fig. 3. Model performance over original and adversarial images.

TABLE II. Model Performance Metrics

| Metric | Value |
|---|---|
| Accuracy | 0.95 |
| Precision | 0.96 |
| Recall | 0.94 |
| F1 Score | 0.95 |
| Training Time (s) | 120.00 |
| Inference Time (s) | 0.50 |
| Model Size (MB) | 1.50 |

Our investigation into masking patterns revealed valuable insights into the benefits of employing random masking during training. The use of random masking, where a subset of input features is randomly masked or set to zero in each training epoch, emerged as particularly advantageous. This approach promotes generalization by preventing the model from overfitting to specific features. Over-reliance on particular features could lead to decreased adaptability and performance degradation when faced with unseen or perturbed data. The adoption of random masking strategies, therefore, contributes to a more robust and versatile model.

A noteworthy observation pertained to a slight increase in training times resulting from the incorporation of feature masking. The additional step of applying masks during each training epoch introduced a minor overhead. However, it is essential to highlight that this increase did not translate into a significant rise in computational resource requirements. The practical feasibility of implementing feature masking in neural network training is underscored by the manageable impact on training times. This finding suggests that the benefits gained in terms of enhanced robustness justify the marginal increase in training duration.

Our study not only highlighted the critical trade-off between accuracy and robustness associated with varying masking ratios but also emphasized the advantages of employing random masking patterns to foster model generalization. Furthermore, the observed increase in training times, while present, did not pose a significant obstacle to the practical implementation of feature masking in neural network training, thereby affirming its feasibility for real-world applications.

## VIII. Results and Discussion

### A. Trade-off Between Accuracy and Robustness

The trade-off observed between accuracy and robustness in adversarial defense strategies can be expressed mathematically.

Fig. 4. Precision, recall and F1 score over epochs.



Fig. 5. Confusion matrix.

Let $\text{Acc}_{\text{baseline}}$ represent the accuracy of the baseline model, $\text{Rob}_{\text{baseline}}$ denote its robustness, and $\text{Mask}_{\text{ratio}}$ be the masking ratio. The relationship can be formalized as follows:

$$\text{Rob}_{\text{masked}} = f(\text{Mask}_{\text{ratio}}, \text{Acc}_{\text{baseline}}) \qquad (24)$$

Here, $f$ is a function that captures the complex interplay between the masking ratio and the baseline accuracy in determining the robustness of the masked model against adversarial attacks. This mathematical representation underscores the necessity of carefully choosing the masking ratio to achieve an optimal balance.

### B. Impact of Feature Masking on Generalization

The study suggests that feature masking, especially with random patterns, favors model generalization but may lead to reduced performance on conventional benchmarks. This can be represented mathematically using the concept of regularization. Let $\mathcal{L}_{\text{masked}}$ denote the loss function for the masked model, and $\lambda$ represent a regularization parameter:

$$\mathcal{L}_{\text{masked}} = \mathcal{L}_{\text{baseline}} + \lambda \cdot \text{Reg}_{\text{masked}} \qquad (25)$$

Here, $\mathcal{L}_{\text{baseline}}$ is the loss of the baseline model, and $\text{Reg}_{\text{masked}}$ represents the regularization term induced by the feature masking. The addition of the regularization term encourages the model to generalize well beyond the training data, but the choice of $\lambda$ becomes crucial in balancing this regularization against benchmark performance.

### C. Avenues for Future Research

The suggestion of exploring the combination of feature masking with other defense mechanisms implies a potential synergy in adversarial defense strategies. Let $\text{Def}_{\text{combined}}$ represent the effectiveness of the combined defense mechanisms, and $\text{Def}_{\text{mask}}$ and $\text{Def}_{\text{other}}$ denote the effectiveness of feature masking and the other defense mechanism individually:

$$\text{Def}_{\text{combined}} = g(\text{Def}_{\text{mask}}, \text{Def}_{\text{other}}) \qquad (26)$$

The function $g$ encapsulates the synergistic effects and interactions between different defense mechanisms, highlighting the need for further exploration in this domain.

### D. Generalizability Across Datasets and Architectures

While the experiments focused on the MNIST dataset, the generalizability of findings to other datasets and model architectures can be expressed mathematically. Let $\text{Gen}_{\text{dataset}}$ represent the generalizability to a specific dataset, and $\text{Gen}_{\text{architecture}}$ denote the generalizability to a particular model architecture:

$$\text{Gen}_{\text{combined}} = h(\text{Gen}_{\text{dataset}}, \text{Gen}_{\text{architecture}}) \qquad (27)$$

The function $h$ captures the combined effect of dataset characteristics and model architecture on the generalizability of the findings.

The study provides valuable mathematical insights into the interplay of key factors in adversarial defense strategies. These formulations help articulate the trade-offs, implications, and potential synergies in a quantitative manner, paving the way for more rigorous analysis and future research directions in the field of adversarial machine learning.

### E. Effect of Masking Ratio

The observed decline in both training and test accuracy with an increase in masking ratio (10% to 50%) can be mathematically represented. Let $\text{Acc}_\text{train}$ and $\text{Acc}_\text{test}$ represent the training and test accuracy, respectively, and $\text{Mask}_\text{ratio}$ denote the masking ratio. The relationship can be expressed as:

$$\text{Acc}_\text{train/test} = g(\text{Mask}_\text{ratio}) \tag{28}$$

where $g$ is a function capturing the impact of the masking ratio on accuracy. Additionally, the improvement in robustness against adversarial attacks ($\text{Accuracy}_\text{under attack}$) with increasing masking ratio reflects the trade-off:

$$\text{Accuracy}_\text{under attack} = h(\text{Mask}_\text{ratio}, \text{Acc}_\text{baseline}) \tag{29}$$

Here, $h$ encapsulates the relationship between the masking ratio, baseline accuracy, and the model's robustness under adversarial attacks.

### F. Random vs. Fixed Masking

Comparing random and fixed masking patterns involves analyzing their impact on adversarial robustness. Let $\text{Accuracy}_\text{random}$ and $\text{Accuracy}_\text{fixed}$ denote the accuracy under attack for random and fixed masking, respectively, both at a 30% ratio. The relationship can be expressed as:

$$\text{Accuracy}_\text{random} = f(\text{Mask}_\text{ratio}, \text{Pattern}_\text{random}) \tag{30}$$

$$\text{Accuracy}_\text{fixed} = f(\text{Mask}_\text{ratio}, \text{Pattern}_\text{fixed}) \tag{31}$$

Here, $f$ captures the influence of masking ratio and specific masking patterns on adversarial robustness. The superiority of random masking suggests its effectiveness in preventing the model from overfitting to fixed unmasked features.

### G. Baseline Comparison

The vulnerability of the baseline model to adversarial attacks, despite exhibiting high accuracy in normal conditions, can be expressed as:

$$\text{Robustness}_\text{baseline} = 1 - \text{Accuracy}_\text{under attack, baseline} \tag{32}$$

This highlights the significance of defensive strategies like feature masking in enhancing the model's robustness in scenarios where adversarial attacks pose a threat.

### H. Training Time Increase

The increase in training time with higher masking ratios can be quantified. Let $\text{Time}_\text{baseline}$ denote the training time for the baseline model, and $\text{Time}_\text{masked}$ represent the training time for the masked model. The relationship can be expressed as:

$$\text{Time}_\text{masked} = i(\text{Mask}_\text{ratio}, \text{Time}_\text{baseline}) \tag{33}$$

Here, $i$ captures the impact of the masking ratio on training time. The more pronounced increase with random masking suggests the additional computational overhead associated with its dynamic nature.

### I. Choosing Masking Ratio

The selection of the appropriate masking ratio involves a trade-off between accuracy and robustness. Let $\text{Utility}_\text{application}$ represent the utility for a specific application, combining accuracy and robustness requirements:

$$\text{Utility}_\text{application} = j(\text{Acc}_\text{test}, \text{Accuracy}_\text{under attack}) \tag{34}$$

Here, $j$ is a function that encapsulates the application-specific requirements, guiding the choice of the optimal masking ratio.

### J. Potential for Further Research

The results indicating the potential for further research can be framed mathematically. Let $\text{Potential}_\text{research}$ represent the potential for further research, considering more sophisticated masking strategies ($\text{Mask}_\text{sophisticated}$), combining feature masking with other defense techniques ($\text{Def}_\text{combined}$), and extending the approach to more complex datasets and models.

A quantitative understanding of the observed effects, trade-offs, and potential for further research in the context of feature masking and adversarial defense strategies.

## IX. CONCLUSION

In summary, our investigation into fortifying adversarial defense in neural networks through the amalgamation of feature masking and gradient manipulation, with a focus on the MNIST dataset, has provided noteworthy insights. The primary aim was to evaluate the efficacy of this approach in enhancing the model's resilience against adversarial attacks, a critical concern in the realm of AI security.

The baseline model, devoid of feature masking, exhibited a commendable accuracy of 98% on the MNIST test set. However, its susceptibility to adversarial attacks was starkly apparent, evidenced by a substantial performance decline to 60% accuracy under Fast Gradient Sign Method (FGSM) attacks. Conversely, models incorporating feature masking displayed varying levels of improved robustness:

10% masking: Despite a marginal decrease in test accuracy to 96%, the model showcased enhanced resilience, maintaining a 75% accuracy under adversarial conditions.

30% masking: A slight dip in test accuracy to 94% was observed, but the model exhibited further improvement in robustness, achieving 80% accuracy against adversarial attacks.

50% masking: While this level led to a more significant accuracy reduction to 92%, it offered the highest defense, reaching an 85% accuracy against attacks.

Additionally, the introduction of feature masking introduced a crucial trade-off between standard accuracy and adversarial robustness. This trade-off holds pivotal significance in applications where the reliability and security of AI models are paramount, such as in autonomous systems and the healthcare industry.

The utilization of a random masking pattern uncovered potential benefits in improving model generalization and resistance against adversarial manipulation. Looking forward, the research holds expansive and promising future prospects. These include exploring advanced feature masking techniques, potentially adaptive or dynamic in nature, and integrating feature masking with other adversarial defense strategies like adversarial training. Moreover, extending the methodology to more intricate datasets and deepening our comprehension of adversarial vulnerabilities in neural networks represent critical strides.

Ultimately, the practical application of these findings in real-world scenarios, especially in high-stakes fields, would signify a substantial advancement in the realms of AI and machine learning.

## REFERENCES

[1] A. Madry, A. Makel, L. Tsipras, and A. Vladu, "Adversarial logit pairing for detecting adversarial examples," *arXiv preprint arXiv:1705.07201*, 2017.

[2] Z. Zhao, S. Gong, and Z. Wang, "Meshed-MD: Security verification and defense for deep neural networks using meshed tensorflow," *arXiv preprint arXiv:2006.08610*, 2020.

[3] E. Wong and J. Z. Kolter, "Distilling robustness in deep neural networks," *arXiv preprint arXiv:1804.03538*, 2018.

[4] A. Madry, A. Makel, L. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.

[5] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Improving the robustness of deep neural networks to adversarial examples," *arXiv preprint arXiv:1412.6574*, 2014.

[6] A. Barz, et al., "Spectral Signature Mismatch for Adversarial Example Detection," *arXiv preprint arXiv:1901.00534*, 2019.

[7] J. H. Metzen, et al., "Input Gradient Consistency for Adversarial Example Detection," *arXiv preprint arXiv:2001.03405*, 2020.

[8] D. Xu, et al., "Adversarial Example Detection via Unsupervised Learning," *arXiv preprint arXiv:2006.06684*, 2020.

[9] J. Alayrac, et al., "Defending Deep Learning with Defensive Adversarial Training," *arXiv preprint arXiv:1804.01444*, 2018.

[10] Y. Zhao, et al., "GAN-Based Gradient Matching for Improved Adversarial Training," *arXiv preprint arXiv:1906.01840*, 2019.

[11] M. Henaff, et al., "Meta-Learning Adversarial Training," *arXiv preprint arXiv:1903.10384*, 2019.

[12] T. Tram, et al., "Ensemble Adversarial Training: Threats and Defenses," *arXiv preprint arXiv:1906.11640*, 2019.

[13] Y. Zhao, et al., "GAN-Based Gradient Matching for Improved Adversarial Training," *arXiv preprint arXiv:1906.01840*, 2019.

[14] T. Tram, et al., "Ensemble Adversarial Training: Threats and Defenses," *arXiv preprint arXiv:1906.11640*, 2019.

[15] S. M. Moosavi-Dezfooli, et al., "DeepFool: A Model-Independent Adversarial Attack on Deep Learning," *arXiv preprint arXiv:1511.04599*, 2016.

[16] S. M. Moosavi-Dezfooli, et al., "DeepFool: A Model-Independent Adversarial Attack on Deep Learning," *arXiv preprint arXiv:1511.04599*, 2016.

[17] Y. Liu, et al., "FDA-based Adversarial Example Detection with Explainability," *arXiv preprint arXiv:2307.03765*, 2023.

[18] X. Wang, et al., "ODAE: Anomaly Detection for Adversarial Examples Using Autoencoders," *arXiv preprint arXiv:2306.01420*, 2023.

[19] Z. Liu, et al., "Explainable Gradient Consistency for Adversarial Example Detection," *arXiv preprint arXiv:2308.03185*, 2023.

[20] H. Zhang, et al., "Adversarial Weight Pruning for Robust Deep Learning," *arXiv preprint arXiv:2305.08465*, 2023.

[21] H. Jiang, et al., "Transferable Adversarial Training via Cross-domain Knowledge Distillation," *arXiv preprint arXiv:2301.10058*, 2023.

[22] Y. Zhao, et al., "DEAD: Dynamic Ensembles for Adversarial Defense," *arXiv preprint arXiv:2307.07056*, 2023.

[23] B. Xu, et al., "Carlini and Wagner Attack Mitigation via Feature Disentanglement," *arXiv preprint arXiv:2309.04057*, 2023.

[24] Y. Liu, et al., "Adaptive Smoothing for DeepFool Defense," *arXiv preprint arXiv:2305.02056*, 2023.

[25] B. Chen, et al., "DeepFool Defense through Multi-Scale Gradient Filtering," *arXiv preprint arXiv:2307.03562*, 2023.

[26] Y. Zhao, et al., "Towards Certified Robustness against Norm-Bounded Adversarial Examples," *arXiv preprint arXiv:2309.03980*, 2023.

[27] Z. Wang, et al., "Feature Pruning for Efficient Adversarial Training," *arXiv preprint arXiv:2306.03761*, 2023.

[28] Z. Liu, et al., "Uncertainty-Aware Adversarial Defense via Ensemble Diversification," *arXiv preprint arXiv:2305.09654*, 2023.

[29] R. Gu, et al., "Adversarial Training with Wasserstein Distance Divergence," *arXiv preprint arXiv:2307.07178*, 2023.

[30] Y. Zhao, et al., "Towards Certified Robustness against Norm-Bounded Adversarial Examples," *arXiv preprint arXiv:2309.03980*, 2023.

[31] Z. Wang, et al., "Feature Pruning for Efficient Adversarial Training," *arXiv preprint arXiv:2306.03761*, 2023.

[32] Z. Liu, et al., "Uncertainty-Aware Adversarial Defense via Ensemble Diversification," *arXiv preprint arXiv:2305.09654*, 2023.

[33] R. Gu, et al., "Adversarial Training with Wasserstein Distance Divergence," *arXiv preprint arXiv:2307.07178*, 2023.

[34] P. Battaglia, J. B. Hamrick, V. Bapst, A. Fontaine, D. Sanchez-Gonzalez, K. Brown, ... and A. Santoro (2018). Relational inductive biases for deep learning. *arXiv preprint arXiv:1802.00302*.

[35] P. Gärdenfors, and R. Pfeifer (2000). Conceptual artifacts in animal minds: A neuro-symbolic approach. *Cognitive Science*, 24(4), 3.

[36] Cubuk,RandAugment: A Simple Data Augmentation Method for Deep Learning,2019

[37] Ganesh Ingle and Sanjesh Pawale, "Generate Adversarial Attack on Graph Neural Network using K-Means Clustering and Class Activation Mapping" International Journal of Advanced Computer Science and Applications(IJACSA), 14(11), 2023.

[38] Ingle, G.B., Kulkarni, M.V. (2021). Adversarial Deep Learning Attacks—A Review. In: Kaiser, M.S., Xie, J., Rathore, V.S. (eds) Information and Communication Technology for Competitive Strategies (ICTCS 2020). Lecture Notes in Networks and Systems, vol 190. Springer, Singapore.

# Automated Paper-based Multiple Choice Scoring Framework using Fast Object Detection Algorithm

Pham Doan Tinh, Ta Quang Minh
School of Electrical and Electronic Engineering
Hanoi University of Science and Technology
Hanoi, Vietnam

*Abstract*—**Optical mark reader (OMR) technology is an important research topic in artificial intelligence, with a wide range of applications such as text processing, document recognition, surveying, statistics, and process automation. Researchers have proposed many methods employing either traditional image processing and statistics or complex machine learning models. This paper presents a feasible solution for the OMR problem. It uses a fast object detection model to detect markers effectively and then segment the answer sheet into smaller regions for the mark reader model to recognize the user's selections accurately. The experimental results on actual answer sheets from college exams show that the error is less than 0.5 percent, and the processing speed can achieve up to 50 answer sheets per minute on standard core i5 personal computers.**

*Keywords*—*Optical mark reader; multiple choice exam; automatic scoring; segmentation; fast object detection*

## I. Introduction

In the era of digitalization and automation, the education sector has attracted significant attention due to its potential to revolutionize traditional educational methods by incorporating cutting-edge technologies to improve the quality of education and academic management. In teaching, evaluating learning progress and the assessment of the learners is very important. Automating this process by applying technology such as an Optical Mark Reader (OMR) attracts the attention of many researchers and organizations. OMR technology has become essential for the automatic multiple choice scoring system, especially in large-scale competitions.

OMR is now widely used for exams or surveys with multiple choice answers [1]. According to Zhang et al. [2], this is the most common type of exercise used in education. This technology focuses on rapidly detecting data extracted from filled-in forms created with a pencil or pen. OMR technology involves the use of Multiple Choice Questions (MCQs) in exams, which allows quick results for students, serves as a tool for teachers and educational institutions to apply in their exams, reduces the need for manual labor, and improves performance. OMR initially appeared as a dedicated hardware solution [3, 4, 5, 6] or using paid resources [7, 8, 9, 10]. These approaches have often been studied before. But then, software solutions [11, 12, 13, 14, 15] appeared along with the development of technology, gradually replacing specialized hardware devices. OMR approaches can be divided into two main categories: Using conventional image processing [4, 16, 17, 18], and using artificial intelligent machine learning [11, 19, 20]. In conventional image processing approaches, first, they adjust the orientation of the input image [5, 11, 20], then apply the segmentation techniques to search for areas that need identification [5, 14]. After that, they detect whether the answer area is circled based on the grayscale level [21, 22] or the number of pixels in the area [18, 23, 24]. These approaches are easy to build and have a short implementation and runtime. However, they may need to fully capture the complex attributes and variations of each specific test, leading to low accuracy.

The limitation of conventional methods has led to growing interest in deep learning methods, especially the convolutional neural networks (CNN), which have demonstrated superior image processing and recognition capabilities. Deep learning offers the potential for many research fields such as image and signal processing [25, 26]. It provides more accurate and robust OMR systems capable of handling diverse types of tests. In addition to the processing algorithms used in the pure image processing approach, this method builds a neural network suitable for the problem. The classification techniques [13, 19, 23, 27] are commonly used. This technique can accurately and quickly identify an answer box to identify whether an answer is selected.

In addition, the input images may come from many sources like cameras, webcams [2, 28] or from smartphones [16, 19, 20, 29]; this factor also dramatically affects construction costs and model implementation time. Models using many image formats and sources will save time and effort and reach more users.

Several methods of deploying the system into software [4, 12] on desktop or mobile devices have built a relatively complete system. The benefit of this is that it can be used flexibly in many places and has high practical applications. These systems often require users to print or create exam papers using predefined software [9]. However, this must ensure excellent and stable performance because it is difficult to maintain, modify, and add features.

According to Sumit Tiwari and colleagues [30], manipulation of OMR board data is shared and affects exams nowadays. This form of data tampering has not been taken into account by existing systems. This article aims to use an algorithm to encode the characteristics of the answers and information students have highlighted in the answer sheet. Then, create a QR code and use that QR code to evaluate whether the exam paper is fake or not. A novel method that achieves successful research results can be applied in practice.

This article addresses the above research limitations by proposing a deep learning method based on the YOLO (You Only Look Once) algorithm to score multiple choice tests

accurately. We use YOLOv8 because it stands out as the fastest model with lower parameters compared to the other versions [31]. This study uses a data set of real-life multiple choice test sets and training and testing processes to create a powerful and effective model. The contributions of the research include flexible use of input images, low implementation costs, and high accuracy requirements, which is important to propose a fast, easy-to-use method with the use of an optimal resource.

The rest of the article is presented as follows: Section II offers the proposed system architecture and detailed algorithm implementation. Section III presents the test and the results evaluations. The final section concludes the article with the future direction of the system given in Section IV.

## II. Method

### A. Answer-sheet design

The answer sheets given to students in each exam are designed as shown in Fig. 1. The above form was redesigned from the answer sheet in Vietnam's national high school exam.
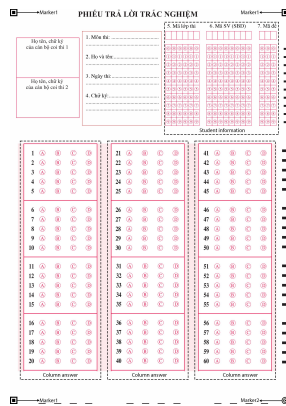


Fig. 1. Sample answer sheet.

The answer sheet has the student's registration number, exam code, and exam-class sections to get information about students and exam questions, making the process of statistics and data processing easier when the data set is large. In addition, the answer section has a maximum of 60 questions, and you can optionally score specific questions, which is suitable for multiple-choice exams.

### B. Overall System Implementation

In Fig. 2, we build a general system diagram for the system analyzed in the previous section. This diagram can be used to understand how the system's components interact with each other, as well as provide an overview of the system's architecture and functionality.

We propose to divide the method into phases: Segment and preprocess data phase, Labeling phase, Training phase, and Online Recognition phase. These stages are presented in detail in the following sections.

### C. Segmentation and Pre-processing Phases

Before entering recognition, to achieve a balanced accuracy and performance, YOLO recommends that the model's input

image be sized 640x640. The model will even resize the input image to have the most significant side size set to 640 and maintain the original aspect ratio. Because of this, if the image is not segmented into small parts, small details will be lost when the input image is trained. Therefore, we improved the accuracy by focusing on the desired portion of the original input image (segmenting the portion and keep the original resolution).

We can find the constant lines that surround the blocks and shapes to segment the input image into student's information section and the question answer choice section as showed in Fig. 3:



Fig. 3. Segmented image.

Segmented image components after cropping will be re-sized prior to recognition. Here, we will resize to have the most significant edge size set to 640 and keep the same proportions.

The training answer sheets are divided into two sets: training and validation with 85% and 15%, respectively.

### D. Labeling phase

After having preprocessed data, we build a labeling process for the model. Labeling is defining bounding boxes around class types in the image and placing captions for each box. The model can be trained to detect and classify classes in the following training phase by accurately locating the courses in the photo. Here, we label each cropped image with the LabelImg software.

*1) Marker:* In reality, the input images can be skewed, rotated, etc. The coordinates of the markers that help us specifically handle these problems will be presented in the next part. We placed three markers in three corners: Top left, top right, and bottom left of the exam paper with the same shape and labeled them as "marker1". The marker in the lower right corner using other shape type and labeled as 'marker2' shown in Fig. 4.

*2) Question-Answer Section:* In this section, each question will have four answer options; each question can have many correct answers, so we will have $2^4$ cases where the answer is selected. Therefore, we use 16 labels, encoded in bits 0 and 1, shown in Table I and labeled as in Fig. 5.

*3) Student Information Section:* Student information includes the exam class code, student registration number, and exam code. These fields are identified by integer numbers from 0 to 9. Therefore, we use 10 labels, shown in Table II and labeled as in Fig. 6.

### E. Training phase

During the training process, there are several main steps as the following description:

Fig. 2. System overview diagram.



Fig. 4. Marker-labeled image.

TABLE I. LABEL ANSWER

| Label | Value | Label | Value |
|-------|-------|-------|-------|
| 1000 | A | 0101 | B and D |
| 0100 | B | 0011 | C and D |
| 0010 | C | 1110 | A, B and C |
| 0001 | D | 1101 | A, B and D |
| 1100 | A and B | 1011 | A, C and D |
| 1010 | A and C | 0111 | B, C and D |
| 1001 | A and D | 1111 | A, B, C and D |
| 0110 | B and C | 0000 | Not selected |



Fig. 6. Labeling student information,



Fig. 5. Labeling answers.

model. Metrics such as loss, mAP are monitored to assess the model. These parameters are explicitly described in Table III.

TABLE III. CUSTOM TRAINING MODEL

| Parameter | Value |
|-----------|-------|
| Model | YOLO |
| Image size | $640 \times 640$ |
| Number of Epochs Trained | 150 |
| Batch Size | 16 |
| Number of classes | 29 |

   *a) How to use the training set:* The training set accounts for 85% of the total data collected. This data set includes labeled images that correspond to each class.

   *b) Select parameters:* YOLOv8's training configuration contains parameters such as Number of classes, Image size, Number of epochs, and Batch size. The training process is monitored to evaluate the progress and effectiveness of the

TABLE II. LABEL INFO

| Label | Value |
|-------|-------|
| 0 | 0 |
| 1 | 1 |
| 2 | 2 |
| 3 | 3 |
| 4 | 4 |
| 5 | 5 |
| 6 | 6 |
| 7 | 7 |
| 8 | 8 |
| 9 | 9 |
| unchoice | Not selected |

*c) Model Optimization:* After completing model training, the next step is to test and fine-tune the model. This step is done to ensure that the model can perform well in new tests and is accurate in different types of questions.

The performance and accuracy of the model are evaluated through commonly used parameters in Machine Learning in general and object recognition problems in general: Precision, recall, mean precision (AP) and mean average precision (mAP) [32]:

$$Precision = \frac{TP}{TP + FP} \qquad Recall = \frac{TP}{TP + FN} \qquad (1)$$

In the Formula 1:

- $TP$ : Number of cases correctly predicted as Positive.

- $FP$ : Number of cases predicted to be Positive but actually Negative.

- $FN$ : Number of cases predicted to be Negative but actually Positive.

From Precision and Recall, we calculate the average accuracy of the object detection model:

$$AP = \sum_{i=0}^{i=n-1} [Recalls(i) - Recalls(i+1)] \times Precisions(i) \qquad (2)$$

Thence inferred:

$$mAP = \frac{1}{h} \sum_{i=1}^{i=n} AP_k \qquad (3)$$

In the Formulas 2 and 3:

- $h$ is the number of classes

- $Recalls(i)$ and $Precisions(i)$ are the value of the $i^{th}$ element of the Recalls and Precisions array

- $AP_k$ is the AP value of the $i^{th}$ class

#### F. Online Recognition Phase

The input image is taken directly from the camera or smartphone, so it is impossible to avoid the cases where the input image has different angles and distances from the camera to the answer sheet or cases where the input image is blurry,

incorrect and misaligned. This stage's purpose is to process the image to extract the part of the image that only contains multiple-choice answer sheets. The test paper must be aligned in the most appropriate direction, brightness, and color to be included in the identification model.

First, predict the input image, the target to identify four markers, and we get position marker. After the YOLO model recognized four markers (3 square markers and one circle marker), we got the coordinates of the four markers on the original exam paper. Note that the order of the detected angles is unconventional. Because of the above reason, we need to rearrange the four corners in the correct order. Based on the Position Marker (PM), we determine the direction of the image by placing three markers1 in positions: top-left, top-right, bottom-left, and marker 2 in the bottom-right position. Therefore, to retrieve the part of the image that only contains multiple-choice answer sheets, we rotate the image so that marker 2 is always in the bottom-right position.

Call the top left point $P_1$, the top right point $P_2$, the bottom right point $P_3$, the bottom left point $P_4$. Suppose that each point is defined by the coordinates (x,y):

$$PM = \{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)\} \qquad (4)$$



Fig. 7. Illustration of the input image.

Fig. 7 illustrates the first image of the process. After obtaining the position of marker 2 through identification, we consider this position to be the new bottom-right position. Then rotate the remaining corners according to this marker2 position. $P_4'$ is the new bottom-left position. In fact, $P_4'$ can be in many places around $P_3$. We need to determine the rotation angle $\alpha$, wherever $P_4'$ is.

First, determine the coordinates $P_4$ among the 3 marker1 coordinates. Because $P_4$ is considered a bottom-left point, based on the distance, we determine $P_4$ is the point with the shortest distance to $P_3$, specifically $d_1 < d_3 < d_2$: $P_4 = \{(x_i, y_i) \subset PM | d_{P4} = min\{d_1, d_2, d_3\}\}$. In the next step, determine the coordinates $P_4'$, $P_4'$ at the new bottom-left position, with a distance equal to $P_4$ to $P_3$, so the coordinates of $P_4'$ are always equal to:

$$\begin{cases} x_{P_4'} = x_{P_1} - d_1 \\ y_{P_4'} = y_{P_3} \end{cases} \qquad (5)$$

From the coordinates $P_4'$, applying the trigonometric formula for triangle $P_4'P_4P_3$, we calculate angle $alpha$ with $d_{P_4P_4'}$ is the distance from $P_4$ to $P_4'$:

$$\alpha = cos^{-1} . \frac{2d_1^2 - d_{P_4P_4'}^2}{2d_1^2} \qquad (6)$$

Rotating the input image with angle $\alpha$, we obtain a new image rotated in the correct direction. After corner points have been identified and target points have been calculated, we use the getPerspectiveTransform and warpPerspective functions to transform and align the input image. This image processing step, which includes markers, produces a transformed image that only contains the answer sheet. Additionally, the image is correctly rotated.

We will obtain an image that only includes the answer sheet and has been transformed accurately. We tested our model on a variety of image samples, including different orientations and lighting conditions, to obtain an image that only contains the answer sheet. The model still works well in most cases. The algorithm used in this procedure is shown in Algorithm 1.

---

**Algorithm 1** Image Preprocessing

---

**INPUT**: Input image
**OUTPUT**: Preprocessed image
1: **begin**
2: Read $input\_image$
3: Recognition $input\_image$
4: Position marker
5: PM = Initial array position marker
6: **if** (PM has 3 marker1) and (PM has 1 marker2) **then**
7:       Find $d_{P_4} = min\{d_1, d_2, d_3\}$
8:       Find $\alpha = cos^{-1} . \frac{2d_{P_4}^2 - d_{P_4P_4'}^2}{2d_{P_4}^2}$
9:       Rotate $input\_image$ with angle $\alpha$
10:      Find destination corners (DC):
    DC = [ [0, 0], [ maxWidth, 0], [maxWidth, maxHeight], [0, maxHeight] ]
11:      Get the background removed image from the DC: $image\_extracted$. Using **getPerspectiveTransform** and **warpPerspective**
12:      return $image\_preprocessed = image\_extracted$
13: **else**:
14:     $break$
15: **end**

---

Fig. 8 shows the final result of Algorithm 1. With this processed image, the following segmentation and identification of each component is much easier.

Get the image after preprocessing, crop the image to get: column answer image and student information image, recognition these images. Based on recognition results, we can extract the information of students and the answers from each answer sheet of the exam. Subsequently, comparisons with the correct answers associated with each exam class code were made, allowing each candidate to receive an automated scoring process. In addition, a threshold coefficient called $\theta$ was introduced. This threshold is the decisive parameter that determines the confidence level required to consider a prediction to be



Fig. 8. Input image and pre-processed image.

"correct". If the confidence exceeds the specified threshold, the result of the recognition will be confirmed as accurate. On the contrary, if the confidence falls below the threshold, the prediction is considered wrong. The algorithm below describes the systematic identification and scoring process:

---

**Algorithm 2** Recognition and Grading

---

**INPUT**: Pre-processed image
**OUTPUT**: Recognized images and mark
1: **begin**
2: Read $preprocessed\_image$
3: Segmentation $preprocessed\_image$: $info\_student$ and $column\_answer$
4: Recognition $info\_student$ and $column\_answer$
5: Threshold = $\theta$
6: **If** confidence $\geq$ threshold **then**
7:     Insert to database
8:     Write mark
9: **else**
10:     Issue a warning
11: **end**:

---

### III. EXPERIMENTAL RESULTS

A set of experiments was performed to evaluate the accuracy of the methods presented and the automated scoring systems. With an algorithm written in Python, the automatic scoring system is evaluated on a laptop running on an Intel Core i5 $11^{th}$ processor with 16GB RAM, recognizing input images including three parts that need to be recognized: marker, information, student, and answer information. The system will take a variable amount of time for the recognition process proportional to the number of questions on the answer sheet. When identifying votes with fewer answers, it will take less time. We experimented and calculated that the average time to recognize a 60-question answer sheet is 1.2 seconds.

The model is used to predict results for new tests. These results are transmitted to the scoring system to produce the final results shown in Fig. 9:

After testing and refining the model and continuing to train, our team achieved the following results after training in Fig. 10:

The Confusion Matrix chart shows the confusion between classes in the entire system. It can be seen that the confusion model is very little, shown in points other than the main diagonal (representing noise) and mainly confusion. between the

Fig. 9. Precision and recall.



Fig. 10. Confusion matrix normalized chart.

background and the labels, not between the labels themselves, and this confusion level has very low reliability (0.01, 0.02, 0.03...).

According to the chart, we can see that the main diagonal is very thick and almost reaches 1, which means the model has high accuracy because the main diagonal of the matrix represents the number of cases in which the model correctly classifies objects into corresponding classes. The Fig. 11 depicts many graphs of the results after training the process.



Fig. 11. Training results chart.

The $train/box\_loss$ and $val/box\_loss$ plots help us evaluate how well the model locates classes by measuring the difference between the predicted bounding-box coordinates and the class in reality. The decreasing trend in these two charts shows that the model is improving its accuracy in identifying the correct location of labels in the training and validation data sets.

The $train/cls\_loss$ and $val/cls\_loss$ plots illustrate the classification error. This histogram measures the difference between the predicted class probability and the actual label. The decreasing $cls\_loss$ plot shows that the model is improving its ability to identify classes correctly.

The $train/dfl\_loss$ and $val/dfl\_loss$ plots illustrate the imbalance between classes with many labels and classes with few labels. This histogram corrects the difference between the predicted probability and the target probability, especially on data sets with class imbalance. The decreasing $dfl\_loss$ plot shows that the model made more balanced predictions and improved performance.

The precision and recall charts have curves near the highest curve, showing that the model achieves high precision and recall when changing the probability threshold. This indicates that the object recognition model can detect and locate objects.

mAP50 plots: This chart depicts the average accuracy; the model achieves a high value (approaching 1), showing that the model achieves high accuracy in object recognition; mAP50-95 chart If the model reaches a high value, it shows that the model can recognize objects well on many different levels of probability threshold. By monitoring these graphs and continuously improving the model based on the insights we gain, we can train the YOLOv8 model to accurately detect, locate, and classify classes for multiple choice exams.

The test set contains new tests that are not used to train the model. Based on Formula 3 and the results after training, the system can evaluate the performance of a model in a new test and calculate the parameters as described in Table IV:

TABLE IV. MEAN AVERAGE PRECISION

| Parameter | Value |
|---|---|
| Evaluation Metric | mAP |
| Best Metric Scores | 0.996 |

After completing the training phase, analyzing the input image and applying the recognition process, rectangles will be drawn on points that the model recognizes: around markers, student information and selected sentences (fill). The results obtained are shown in Fig. 12.



Fig. 12. Input image prediction.

To get an accurate experiment, we will use a data set such that input images are from different angles and resolutions, taken from many types of devices with diverse light intensities. In addition, on each of those answer sheets, the number of answers varies, with multiple answer choices. To make the test data set diverse, accurate, and most importantly, "realistic," we

used this system to automatically score 300 real-life multiple-choice tests with the semester's final exam at Hanoi University of Science and Technology.

Choosing the threshold value is an integral part of evaluating model performance. The input image is likely wrong when one of the detected model objects has a confidence level below the Threshold threshold that we have previously chosen. Perform the first rough tuning experiment: Choose $0.6 \leq \theta \leq 0.9$, each increment of 0.05 each (see Fig. 13).



Fig. 13. Relationship between error and threshold $\theta$.

According to the results of the diagram above, perform the second experiment for fine tuning (see Fig. 14): Choose $0.75 \leq \theta \leq 0.85$, each mark is 0.01 apart. The purpose is to choose the most suitable $\theta$ value: Our system shows



Fig. 14. Relationship between error and threshold $\theta$ (fine tuning).

high accuracy when input images are of good quality, such as answer sheets taken in suitable, flat, and transparent lighting conditions. However, the system needs better-quality input images. For example, when the answer sheets are blurry, uneven, or have a lot of extra lines due to uneven scanning, the model needs help identifying the answers and information from the student. Especially when the input image lacks corners, markers are lost, leading to the image preprocessing process being unable to process. The accuracy of the system can be significantly reduced if necessary information is lost or contaminated. Based on the analysis results of Fig. 14, we decided to choose $\theta = \mathbf{0.79}$ - a value large enough to have the lowest probability of errors occurring in the data set.
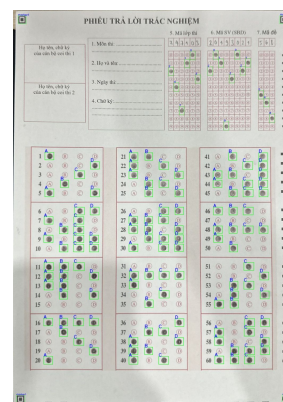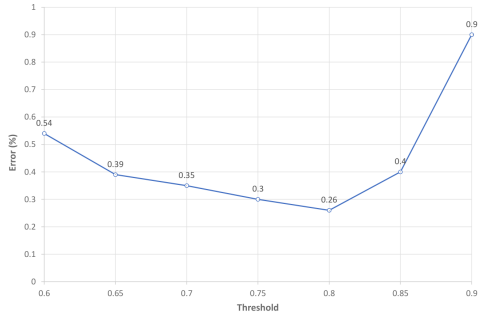
## IV. CONCLUSION

This paper presented an automated paper-based multiple choice grading system with the ulitization of using fast object detection algorithm. Research and experimental results on actual college exams have shown that employing YOLOv8 model together with pre-processing techniques improved the performance of the OMR system with the error rate less than 0.5% and processing time on stand personal computer around 1 second. This can help educational and assessment organizations perform test administration tasks more effectively. The article also highlights challenges and potential development directions. Integrating the system into real-world applications and improving the real-time application ability are also challenges worth considering in the future.

## REFERENCES

[1] E. M. de Elias, P. M. Tasinaffo, and R. H. Jr., "Optical mark recognition: Advances, difficulties, and limitations," *IEEE*, vol. 2, no. 367, 05 July 2021.

[2] L. Zhang, B. Li, Q. Zhang, and I.-H. Hsiao, "Does a distributed practice strategy for multiple choice questions help novices learn programming?" *IEEE*, vol. 15, no. 18, September 2020.

[3] P. Sanguansat, "Robust and low-cost optical mark recognition for automated data entry," *IEEE*, 20 August 2015.

[4] V. Ware, N. Menon, P. Varute, and R. Dhannawat, "Cost effective optical mark recognition software for educational institutions," *IEEE*, vol. 5, no. 2, 21 September 2018.

[5] N. C. D. Kumar, K. V. Suresh, and R. Dinesh, "Automated parameterless optical mark recognition," *IEEE*, vol. 43, pp. 185 – 195, 05 November 2018.

[6] R. H. Hasan and E. A. Kareem, "An image processing oriented optical mark reader based on modify multi-connect architecture mmca," *INJMTER*, March 2015.

[7] A. H. S. Saad, M. S. Mohamed, and E. H. Hafez, "Coverless image steganography based on optical mark recognition and machine learning," *IEEE*, vol. 9, pp. 16 522 – 16 531, 11 January 2021.

[8] A. B. Talib, N. B. Ahmad, and W. Tahar, "Omr form inspection by web camera using shape-based matching approach," *IJRES*, vol. 3, no. 3, March. 2015.

[9] R. S, K. Atal, and A. Arora, "Cost effective optical mark reader," *IEEE*, vol. 74, no. 2, 31 October 2022.

[10] G. R. I. Rasiq, A. A. Sefat, and M. F. Hasnain, "Mobile based mcq answer sheet analysis and evaluation application," *IEEE*, 16 June 2020.

[11] R. Ahad, R. Toufiq, and S. U. Zaman, "Information retrieval from damage omr sheet by using image processing," *IEEE*, 31 December 2020.

[12] J. L. Pérez-Benedito, E. Q. Aragón, and J. A. A. L. Medic, "Optical mark recognition in student continuous assessment," *IEEE*, vol. 9, no. 4, pp. 133 – 138, November 2014.

[13] M. Afifi and K. F. Hussain, "The achievement of higher flexibility in multiple-choice-based tests using image classification techniques," *IJDAR*, vol. 22, pp. 127 – 142, 19 March 2019.

[14] O. Espitia, A. Paez, Y. Mejia, M. Carrasco, and N. Gonzalez, "Optical mark recognition based on image processing techniques for the answer sheets of the colombian high-stakes tests," *IEEE*, vol. 1052, pp. 167 – 176, 09 October 2019.

[15] E. M. de Elias, P. M. Tasinaffo, and R. Hirata, "Alignment, scale and skew correction for optical mark recognition documents based," *IEEE*, 21 October 2019.

[16] H. Tjahyadi, S. Lukas, S. Albert, and D. Krisnadi, "Automated scoring of multiple-choice test using template matching technique," *ILCA*, vol. 55, pp. 1–5, 21 September 2018.

[17] D. Chai, "Automated marking of printed multiple choice answer sheets," *IEEE*, 13 February 2017.

[18] S. Hussmann and P. W. Deng, "A high-speed optical mark reader hardware implementation at low cost using programmable logic," *IEEE*, vol. 11, no. 1, pp. 19 – 30,, November 2015.

[19] B. Haskins, "Contrasting classifiers for software-based omr responses," *IEEE*, 17 December 2015.

[20] T. D. Nguyen, Q. H. Manh, P. B. Minh, L. N. Thanh, and T. M. Hoang, "Efficient and reliable camera based multiple-choice test grading system," *IEEE*, 26 September 2011.

[21] R. Patel, S. Sanghavi, D. Gupta, and M. S. Raval, "Checkit - a low cost mobile omr system," *IEEE*, 07 January 2016.

[22] F. de Assis Zampirolli, "Automatic correction of multiple-choice tests using digital cameras and image processing," *IEEE*, vol. 13, April 2013.

[23] O. Gorokhovatskyi, "Neocognitron as a tool for optical marks recognition," *IEEE*, 06 October 2016.

[24] H. Tjahyadi, S. Lukas, S. Albert, and D. Krisnadi, "A novel optical mark recognition technique based on biogeography based optimization," *IJITK*, vol. 5, no. 2, pp. 331–333, December 2012.

[25] P. D. Tinh and B. H. Hoang, "Wifi indoor positioning with genetic and machine learning autonomous war-driving scheme," *IJACSA*, vol. 13, no. 2.

[26] V. Monga, Y. Li, and Y. C. Eldar, "Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing," *IEEE Signal Processing Magazine*, vol. 38, no. 2, pp. 18–44, 2021.

[27] A. AL-Marakeby, "Multi-core processors for camera based omr," *ILCA*, vol. 68, no. 13, April 2013.

[28] H. Atasoy, E. Yildirim, Y. Kutlu, and K. Tohma, "Webcam based real-time robust optical mark recognition," *IEEE*, November 2015.

[29] P. D. Tinh, B. H. Hoang, and N. D. Cuong, "A genetic based indoor positioning algorithm using wi-fi received signal strength and motion data," *IAES IJAI*, vol. 12, no. 1.

[30] S. Tiwari and S. Sahu, "A novel approach for the detection of omr sheet tampering using encrypted qr code," *IEEE*, 07 September 2015.

[31] W. Zhang, "A fruit ripeness detection method using adapted deep learning-based approach," *IJACSA*, vol. 14, no. 9.

[32] H. T. Ngoc, N. N. Vinh, N. T. Nguyen, and L.-D. Quach, "Efficient evaluation of slam methods and integration of human detection with yolo based on multiple optimization in ros2," *IJACSA*, vol. 14, no. 11.

# EpiNet: A Hybrid Machine Learning Model for Epileptic Seizure Prediction using EEG Signals from a 500 Patient Dataset

Oishika Khair Esha, Nasima Begum*, Shaila Rahman
Department of Computer Science and Engineering
University of Asia Pacific, Dhaka, Bangladesh

*Abstract*—The accurate prognosis of epileptic seizures has great significance in enhancing the management of epilepsy, necessitating the creation of robust and precise predictive models. EpiNet, our hybrid machine learning model for EEG signal analysis, incorporates key elements of computer vision and machine learning , positioning it within this advancing technological domain for enhanced seizure prediction accuracy. Hence, this research aims to provide a thorough investigation using the Bonn Electroencephalogram (EEG) signals dataset as an alternative method. The methodology used in this study encompasses the training of five machine learning models, such as Support Vector Machines (SVM), Gaussian Naive Bayes, Gradient Boosting, XGBoost, and LightGBM. Performance criteria, including accuracy, sensitivity, specificity, precision, recall, and F1-score, are extensively used to assess the efficacy of each model. A unique contribution is the development of a hybrid model, integrating predictions from individual models to enhance the overall accuracy of epilepsy identification. Experimental results demonstrate notable success, with the hybrid model achieving an accuracy of 99.81%. Performance matrices for both classes demonstrate the hybrid model's epileptic seizure prediction reliability. Visualizations, including ROC-AUC curves and accuracy curves, provide a nuanced understanding of the models' discriminative abilities and performance improvement with increasing sample size. A comparative analysis with existing studies reaffirms the advancement of our research, positioning it at the forefront of epileptic seizure prediction. This study not only highlights the promising integration of machine learning in medical diagnostics but also emphasises areas for future refinement. The achieved results open avenues for proactive healthcare management and improved patient outcomes.

*Keywords*—*Epilepsy; seizure prediction; computer vision; hybrid model; electroencephalography; bonn dataset; proactive healthcare*

## I. INTRODUCTION

In this study, we introduced EpiNet, a novel hybrid machine learning model, designed to significantly advance epileptic seizure prediction using EEG signals. EpiNet uniquely combines the strengths of various advanced machine learning techniques, resulting in a model that not only outperforms existing single-model systems in accuracy but also addresses critical challenges in seizure prediction such as high variability in EEG signals and the need for reducing false positives. Our model stands out in its ability to integrate complex patterns from a large dataset of 500 patients, providing a more robust and reliable prediction mechanism. The introduction of EpiNet represents a pivotal step forward in epilepsy management, promising to enhance patient care through more precise and proactive strategies. Despite the advancements in medical science, epilepsy is a prevalent cerebral disease affecting a substantial portion of the global population [1]. The primary mode of treatment involves the use of medications; however, a considerable proportion of individuals diagnosed with epilepsy have difficulties effectively managing their medication, resulting in a substantial decrease in their overall state of life. For some individuals, the consideration of respecting portions of the brain becomes a drastic option to eliminate the seizure focus, yet this measure does not guarantee freedom from seizures. In response to the limitations of existing treatments, there has been a burgeoning interest in the development of clinically effective seizure prediction systems. A successful prediction method could offer timely warnings to patients, allowing for preventive measures or interventions such as electrical stimulation, medication release, or cooling of the seizure focus area. Despite early attempts to predict seizures, the scientific and clinical communities have faced persistent challenges, partly attributed to the absence of a detailed definition of the preictal stage, the critical period preceding a seizure. The unpredictable nature of seizures adds a layer of complexity and concern for individuals living with epilepsy.

This research endeavour aims to tackle the aforementioned issues by investigating other approaches that might enhance the accuracy of seizure prediction. In our earlier conference paper [2], we conducted an extensive review of machine learning and deep learning approaches for epilepsy diagnosis, identifying critical research gaps such as dataset limitations, preprocessing challenges, and the need for ensemble methods. Building upon these insights, our current paper addresses these challenges head-on. Furthermore, our study delves into the critical aspects of data, feature selection, and model selection, drawing on the recommendations outlined in our previous work. Notably, we emphasise the application of ensemble learning, an idea proposed in our conference paper, demonstrating its effectiveness in enhancing prediction accuracy and mitigating false alarm rates. This seamless connection highlights the evolutionary progression of our research agenda, from identifying challenges to proposing practical solutions. The aim is to predict upcoming seizures before their occurrence, facilitating prompt management and mitigating associated dangers. The subsequent sections delve into a comprehensive methodology and findings, contributing significantly to the advancement of seizure prediction research.

Fig. 1 provides a visual representation of the basic technique used for the recording of epileptic seizures through EEG.
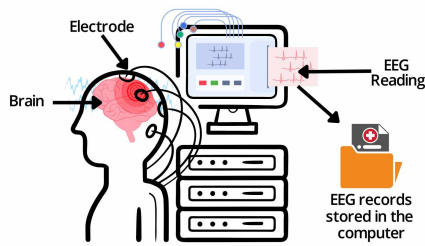
Fig. 1. General methodology for recording epileptic seizures.

During this procedure, electrodes are carefully set on the top of the head in order to assess the electroencephalographic signals, which measure the neural electrical signals produced by the brain. The electrodes are capable of capturing the complex signals produced by neurons in the cerebral cortex.

### A. Key Contributions

In this study, we introduce an innovative machine learning-based approach to epileptic seizure prediction, labelled EpiNet. While significant strides have been made in the realm of epileptic seizure detection, current methods continue to face challenges, and numerous unresolved issues persist. In light of this, our work aims to address a few pivotal questions, outlined below:

*1) Single-Model Limitation in Epilepsy Prediction:* Problem: Existing approaches to epilepsy prediction often rely on single machine learning models, limiting the overall accuracy and robustness of the predictions.

Contribution: Proposed a novel hybrid model that amalgamates predictions from diverse models, demonstrating superior performance with heightened accuracy, sensitivity, and specificity compared to individual models. Despite the limited number of existing hybrid models, our approach stands out by consistently outperforming them, underscoring the effectiveness of our tailored combination of models in epileptic seizure prediction.

*2) Feature Redundancy and Noise:* Problem: The accuracy of prediction models heavily depends on the quality of their features. The presence of redundant or noisy features can compromise the precision of predictions. Identifying and addressing this problem becomes pivotal for enhancing the reliability of epilepsy prediction models.

Contribution: By implementing RFE, the study actively contributes to refining the feature selection process. It goes beyond recognising the issue and presents a strategic solution that decreases noise and improves epilepsy prediction. This work improves prediction model robustness by advancing methodology.

*3) Need for High Accuracy and Reliability:* Problem: For successful patient treatment, medical diagnostics, notably epilepsy prediction, need great accuracy and dependability.

Contribution: Effectively predicted seizures with 99.81% accuracy using EpiNet.

*4) Future Research Directions and Validation:* Problem: The application of current research to a variety of situations and real-world healthcare settings often presents difficulties.

By pointing out the areas that need attention and development, acknowledging these limitations creates the foundation for a significant contribution.

Contribution: This contribution goes above and beyond just identifying limits by providing a clear path for further study. It takes the lead in resolving the issue by outlining concrete measures to improve the generalizability of the model and to verify its efficacy in actual healthcare settings. This prospective strategy sets the research up to be a driving force behind real improvements in the area of epilepsy prediction.

The remainder of this paper is structured as follows: Section II represents the overview of epilepsy. In Section III we have the literature review. The dataset processing and feature extraction is provided in Section IV. The proposed methodology is described in Section V. Section VI analyzes the results of the conducted experiments. Lastly, Section VII concludes the paper with some future works.

## II. OVERVIEW OF EPILEPSY

Epilepsy, an illness that is rather common, has a substantial influence on the lives of millions of people all over the globe. The purpose of this introductory part is to offer a comprehensive viewpoint on the difficulties that epilepsy presents, so laying the groundwork for a more in-depth investigation into the various seizure prediction approaches.

### A. Introduction to Epilepsy

Epilepsy is characterised by recurrent seizures, affecting a considerable portion of the global population. Despite advancements in medical interventions, a substantial number of epilepsy patients face challenges in symptom management with traditional medications. Patients facing resistance to conventional treatments often endure a severely diminished quality of life. Some explore extreme measures, such as brain resection, in pursuit of relief, yet this drastic approach remains ineffective for a notable proportion of individuals. The Fig. 2 visually represents the intricate activity within the brain during an epileptic seizure.



Fig. 2. Illustration of epileptic seizure activity in the brain.

### B. Seizure Types and Symptoms

Epileptic seizures manifest in various types, including focal and generalised seizures, each presenting unique characteristics. A nuanced understanding of these manifestations is essential for accurate diagnosis and prediction. The diverse manifestations of epileptic seizures encompass various types, each characterised by distinct symptoms. Simple Partial Seizures, or Focal Onset Aware Seizures, affect a specific region of the

brain, resulting in altered emotions, sensory perceptions, or movements without loss of consciousness. As depicted in Fig. 3, these seizures may progress to Complex Partial Seizures, or Focal Onset Impaired Awareness Seizures, where consciousness becomes altered, accompanied by involuntary repetitive movements. Notably, these partial seizures can evolve into Generalized Seizures, involving the entire brain.



Fig. 3. Types of seizures.

The visualisation in Fig. 3 offers a thorough and inclusive depiction of the development and unique attributes associated with these many kinds of seizures. This visualisation serves to enhance comprehension and facilitate the categorization of epileptic occurrences. Recognising the diverse symptoms accompanying seizures is critical, emphasizing the need for tailored diagnostic and predictive approaches that consider the individualised nature of epilepsy.

### C. Motivation for Advanced Seizure Prediction Approaches

The primary motivation behind this research stems from the pressing challenges in epilepsy management, particularly the limitations of current diagnostic tools and treatment strategies. One of the key difficulties lies in the unpredictable nature of seizures, which significantly affects patients' quality of life and complicates treatment planning. Current methods lack the precision and foresight needed for effective seizure management. This gap highlights the necessity for more accurate and timely seizure prognosis, where machine learning approaches hold significant promise. Machine learning's ability to analyze complex EEG data and identify patterns indicative of impending seizures presents a transformative opportunity in epilepsy care with minimum costs. Our research with the EpiNet model is directly motivated by these challenges. We aim to leverage the advanced capabilities of machine learning to enhance seizure prediction accuracy, ultimately leading to more personalized and effective epilepsy management strategies. This approach not only addresses a critical need in epilepsy care but also opens up new avenues for research and treatment methodologies in the field.

### III. LITERATURE REVIEW

Kapoor et al. [3] introduce an innovative seizure prediction method using ensemble classifiers and hybrid search optimization, achieving a high accuracy of 96.61% on the CHB-MIT database. Their approach addresses existing limitations, sets a standard for researchers, and explores COVID-19-related data applications for enhanced seizure prediction. Savadkoohi et al.[4] used K-nearest neighbours (KNN) and

support vector machine (SVM) prediction models to predict seizures, demonstrating efficiency, reliability, and flexibility across different frequency ranges. The technique has phase information and directionality issues, despite its merits. [5] introduced a spike rate-based seizure detection approach with great accuracy, improving the quality of life for epileptic patients. However, without a complete deep learning method for prediction, performance may decline.

Usman et al. [6] developed generative adversarial networks using LSTM units to improve sensitivity and reduce false positives. However, anticipation time improvement was inefficient. Author in [7] developed a seizure prediction component using deep learning approaches, improving sensitivity and specificity. The approach has limitations, including a low Signal-to-Noise Ratio (SNR) and dependency on several factors. Emara et al. [8] developed a technique to identify abnormalities in multi-channel EEG signals with high prediction accuracy. This method relied on samples, which was a drawback. Wang et al. [9] pioneered a CNN and DTF-based seizure prediction system, offering potential benefits for epilepsy patients in closed-loop therapy. However, the method was noted for its effectiveness despite time constraints. In a recent study [10], a DWT-transformed EEG data approach, coupled with a DenseNet-LSTM hybrid model, demonstrated improved seizure prediction accuracy, outperforming prior methods on the CHB-MIT scalp EEG dataset. Notably, this integration of Discrete Wavelet Transform and hybrid m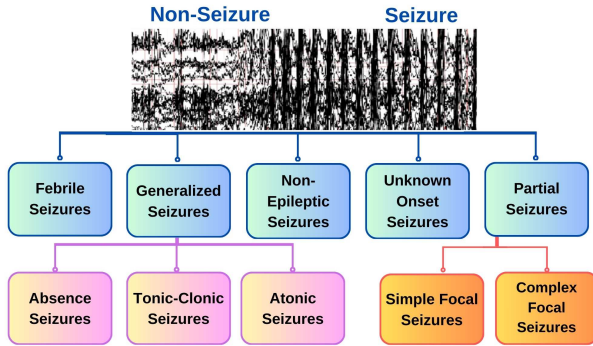odels signifies progress in the field. Viana et al. [11] explored remote subcutaneous EEG monitoring for individualized seizure forecasting.

Behnoush et al. [12] meticulously assessed machine learning algorithms for predicting seizures, emphasizing critical patient identification in emergency department data. A study enhanced signal-to-noise ratio by combining 23 EEG channels into one [13]. Ouichka et al. [14] delved into the challenges of predicting epileptic episodes using iEEG data. Deep learning models, including 3-CNN and 4-CNN, achieved 95% accuracy in autonomous seizure prediction, surpassing previous approaches. In [15], a successful seizure prediction method employs deep learning on preprocessed scalp EEG signals, achieving 92.7% sensitivity and 90.8% specificity in 24 patients. Meanwhile, [16] introduces a novel seizure detection approach using sparse representation with the Stein kernel, leveraging old data to simplify and understand new EEG samples.

The study by [17] employed a two-step strategy, utilizing multi-lead EEG samples to train SE-Net for short-term features and LSTM for long-term characteristics. Adversarial learning enhanced the LSTM feature mapping, resulting in a 5% improvement in classification accuracy on the TUH EEG Seizure Corpus and CHB-MIT databases. This study [18] introduces an innovative approach utilizing SLT and VGG-19 neural networks for precise seizure detection, achieving remarkable 100% accuracy in distinguishing seizure and non-seizure events across seven instances. Notably, its effectiveness extends to improved classification in three- and five-class scenarios, outperforming conventional methods. Tested on the CHB-MIT scalp EEG database, the proposed technique attains a notable 94.3% accuracy in distinguishing seizures from non-seizure episodes.

Studies by [19] highlight LSTM and Random Forest's

efficacy in epileptic episode prediction with 97% and 98% accuracy. [20] conducts a comprehensive literature review emphasizing the critical role of ML in automating epilepsy diagnosis, discussing feature extraction methods, and advocating for relevant characteristics and classifiers. Researchers discovered a novel method for simultaneous epilepsy prediction in adults and children, leveraging a linear mixed model and recording over 1.2 million seizures [21]. Emphasizing the distinct seizure patterns in different age groups, they underscored the need for early diagnosis and treatment to prevent potential brain damage. In a separate investigation, a researcher achieved 100% accuracy in epileptic episode detection using innovative SVM-PCA methodologies, outperforming traditional algorithms [22]. The Local Binary Pattern (LBP) method, assessing key points in EEG data, proves effective for real-time seizure diagnosis [23]. Toraman et al. [24] successfully distinguish preictal and interictal occurrences using SVM, forecasting seizures up to 33 minutes in advance.

Weighted Majority Voting Ensemble (WMVE) stands out by dynamically evaluating and adjusting each classifier, prioritizing accurate categorization, particularly for challenging data. In comparison to Simple Majority Voting Ensemble (SMVE), WMVE demonstrates superior classification accuracy across diverse datasets [25]. Additionally, support vector machines (SVM), particularly when paired with a radial basis function kernel, emerge as highly effective in EEG signal categorization for epilepsy detection [26]. The proposed approach in [27] accurately predicted seizures with an average lead time of 23.6 minutes, utilizing "optimum allocated techniques" for sample selection. In [28], a three-part procedure combining LSTM, regression-based SNR enhancement, and statistical properties achieved superior results (94% accuracy) on the CHB-MIT dataset. Another study [29] employed three machine learning methods for effective seizure detection.

In one approach [30], a phase-space adjacency graph achieved a 97% success rate, while another method using hypergraph analysis reached 93% accuracy. A third method employing deep learning and CNN in phase-space analysis achieved perfect 100% accuracy. Additionally, an ensemble classifier in epilepsy research demonstrated superior performance with a 90% success rate, outperforming others in the range of 85% to 89.5%. Sharma et al. [31] introduce a novel hybrid method, combining higher-order statistics, sensitivity analysis, and the residual wavelet transform, to assess brain signal frequencies affected by transient events. This approach effectively detects and characterizes non-stationary time series alterations in neural activity across different brain areas.

Researchers in [32] developed a machine learning approach for epilepsy prediction using correlation dimension, which converges quickly due to its simplicity. The model predicts seizures by analyzing EEG spike rates, employing a mean filter to smooth spikes and activating an alert when preictal spike numbers exceed a threshold. The suggested technique in [33] predicts seizure activity with 92% accuracy using the CHB-MIT dataset for all patients. Researchers in [34] advocate using reconstructed phase space (RPS) over raw EEG data for seizure detection, citing improved accuracy rates of 95% for tertiary and 98.5% for binary classification. In [35], a two-layer LSTM model achieves an exceptional 98.14% average accuracy, enhancing the quality of life for epilepsy

patients. Additionally, [36] introduces an approach to extract time-frequency features using STFT, addressing CNN and Transformer limitations with an innovative alternating structure for enhanced predictions.

To address the challenge of limited EEG data, the study in [37] employs a deep learning approach, specifically a DCGAN, to generate synthetic EEG data. The proposed method combines transfer learning with popular DL models and utilizes synthetic data for training, showcasing improved epileptic seizure prediction. The study [38] evaluated an EEG-based seizure detection model on CHB-MIT scalp data, demonstrating robustness with sensitivity and specificity values reaching approximately 100%. The updated XG Boost classifier surpassed previous methods, enhancing sensitivity by 0.05% and specificity by 1%. Syed Muhammad Usman et al. [39] created a more sensitive ensemble learning technique for epileptic prediction. Importantly, our technique doesn't use heart rate variability and EEG measurements.

## IV. DATA PROCESSING AND FEATURE EXTRACTION

In this section, we delve into the details of the dataset employed in our study. We'll walk you through the different categories within the dataset, the distinctive features that shape its structure, and the meticulous steps we took to make sure the data is not only accurate but also well-suited for machine learning purposes.

### A. Dataset Description

The Bonn EEG dataset [40], selected for our study, is renowned in epilepsy research for its high-quality and diverse data. This dataset is a benchmark in the field, offering a broad spectrum of EEG signal patterns from various types of epileptic seizures, which is crucial for developing robust and comprehensive seizure prediction models. Its widespread use and proven success in previous studies underscore its reliability and effectiveness. The practicality and accessibility of the Bonn dataset, combined with its ethical soundness, make it an ideal choice for our research, facilitating the advancement of machine learning approaches in epilepsy prognosis. It was collected from the brain activity recordings (EEG) of 500 different individuals. To make the data more manageable for analysis, we divided each 23.6-second EEG recording into 4097 data points. The dataset is neatly organized into 23 segments, each containing a total of 4097 data points. We then mixed things up a bit by shuffling the segments randomly, and within each segment, we have 178 data points. Each of these data points corresponds to a one-second interval in the EEG recording, helping us capture a more detailed snapshot of the brain activity. Consequently, 11,500 rows of data are generated, where each row stands for one second of an individual's EEG recording. A bandpass filter with a frequency range of 0.53–40 Hz and a sampling rate of 173.61 Hz was used in order to do the preprocessing on the dataset first. The preprocessed dataset is made up of electroencephalogram recordings that were taken under a variety of situations, with each condition being assigned a unique class label.

The EEG signals from the Bonn University dataset that represent the five classes are shown in Fig. 4 This picture displays exemplars of the EEG signals that reflect the five
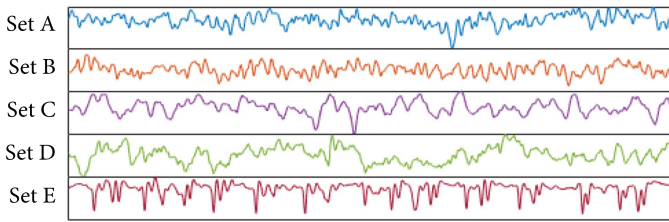
Fig. 4. Data set representing five classes EEG signals.

distinct categories. (a) Set A represents the state of having open eyes, (b) Set B represents the condition of having closed eyes, (c) Set C represents the state of the hippocampal formation during the period between seizures, (d) Set D represents the epileptogenic zone during the period between seizures, and (e) Set E represents the state of having a seizure, also known as the ictal state.

### B. Dataset Structure

The Bonn University Epilepsy dataset comprises five sets, labeled A, B, C, D, and E. Each set contains 100 files, representing recordings from different individuals. Each file corresponds to a single subject and has a duration of 23.6 seconds. The EEG signals in each file are discretized into 4097 data points.

TABLE I. DATASET OVERVIEW

| Class | Patient Status | Setup | Phase |
|---|---|---|---|
| A | Non-Epilepsy | Surface EEG | Open Eyes |
| B | Non-Epilepsy | Surface EEG | Close Eyes |
| C | Epilepsy | Intracranial EEG | Interictal Hippocampal Position |
| D | Epilepsy | Intracranial EEG | Interictal Epileptogenic Zone |
| E | Epilepsy | Intracranial EEG | Ictal |

Table I displays the dataset with category labels. The dataset is categorised into five distinct classes, each representing various situations. As follows:

- Class A: EEG recorded with the patient's eyes open.

- Class B: EEG recorded with the patient's eyes closed.

- Class C: EEG recorded from the healthy brain area where the tumor was not present.

- Class D: EEG recorded from the area where the tumor was located.

- Class E: Seizure activity.

### C. Data Cleaning and Missing Value Handling

For the purpose of ensuring that the data contains no errors, a data cleaning operation was carried out. For the purpose of removing rows from the dataset that were missing values, the dropna() function was used. It was essential to get rid of any data that was erroneous or missing. It made certain that we have trustworthy data.

### D. Outlier Detection and Removal

Using the Local Outlier Factor (LOF) technique, we found and removed any strange data points that may have skewed our estimates. Basically a digital investigator. The "LocalOutlierFactor" algorithm from the scikit-learn module proved to be instrumental in resolving the problem. It helped us find the

data points that didn't seem to fit and get rid of them so they wouldn't mess up the calculations later on.

### E. Feature Scaling and Standardisation

Standardisation of the features was accomplished by using the StandardScaler function that is available in the scikit-learn package. As a component of this, the attributes were normalised by bringing their range to a value of one and with zero serving as the centre of their average. Through the use of a single scale, this step ensured that the features could be compared in a manner that was both accurate and relevant across a variety of characteristic categories.

### F. Feature Selection

Selecting high-quality features is like assembling a winning dish in our machine learning adventure. For the purpose of epilepsy prediction, we are interested in the most relevant ones. Therefore, we used a method known as Recursive Feature Elimination (RFE) to determine the most important characteristics. It's as if we had a handy reference (the RFE class from the scikit-learn package) that highlighted the essential components for precise epilepsy prediction.

*1) Recursive Feature Elimination (RFE):* In the face of information overload, Recursive Feature Elimination (RFE) goes about its business. The process begins with an expansive viewpoint, seeing all characteristics as possible indicators. In order to determine the relative value of each feature, RFE uses a Random Forest classifier as its investigation tool. While RFE iteratively reduces the least significant features, the investigative process starts. Imagine it as a detective selecting the best number of characteristics by sorting through clues and eliminating irrelevant ones. To guarantee that the final collection of features reflects the most important bits of information for the work at hand, RFE's unique technique replicates a seasoned investigator's tactics.

*2) Selected Feature Subset:* When we applied the Recursive Feature Elimination (RFE) method, it helped us identify a subset of the most crucial features. We selected these variables because they demonstrated a remarkable ability to predict epileptic episodes. The idea behind choosing these particular features was to enhance the predictive performance of our models. By focusing on these key variables and narrowing down the feature space, we significantly improved the models' ability to make accurate predictions of epileptic episodes.

## V. PROPOSED METHODOLOGY

In our study, we selected a diverse array of machine learning models, each chosen for its specific strengths in handling the complexities of EEG data. Support Vector Machines (SVM) and Gaussian Naive Bayes were chosen for their effectiveness in high-dimensional spaces and probabilistic approach, respectively. Gradient Boosting, XGBoost, and LightGBM were included for their robustness to overfitting and efficiency in processing large datasets. These models collectively address the challenges of EEG data analysis, such as noise, high dimensionality, and class imbalance. Training and evaluation were challenging due to the intricate nature of EEG signals, with specific strategies employed to mitigate

issues like overfitting and ensure model accuracy and relia-
bility. Figuring out when seizures may take place was the
objective of this study. Therefore, the whole concept hinged on
the utilisation of these ingenious models in order to properly
forecast epileptic episodes.

### A. Model Network

Support Vector Machines (SVM), Gaussian Naive Bayes
(GNB), Gradient Boosting (GB), XGBoost, and LightGBM
are some of the machine learning models that are part of the
application. These models are carefully chosen for their ability
to handle different aspects of the prediction.

### B. Algorithmic Steps

The proposed algorithm comprises the following key steps:

1) Data Loading and Preprocessing:
   - Load the EEG dataset and Preprocess target
     variable for consistency.
   - Drop rows with missing values to maintain
     data integrity.
   - Apply the "Local Outlier Factor (LOF)" algo-
     rithm for outlier identification and exclusion.
2) Feature Scaling and Selection:
   - Scale features using "StandardScaler" for nor-
     malisation.
   - Use Recursive Feature Elimination (RFE)
     with "RandomForestClassifier" to select infor-
     mative features.
3) Dataset Splitting:
   - Extract the dataset into training and test sets
     using "train_test_split" function.
4) Model Training and Evaluation:
   - Train various machine learning models on
     the training set, including SVM, Gaussian
     Naive Bayes, Gradient Boosting, XGBoost,
     and LightGBM.
   - Evaluate trained models using performance
     metrics (accuracy, sensitivity, specificity, pre-
     cision, recall, and F1-score).
5) Hybrid Model Creation:
   - Average predictions from individual models
     to create a hybrid model.
6) Performance Analysis and Visualization:
   - Compare the performance of each model.
   - Plot ROC curves for model performance vi-
     sualization.
   - Generate accuracy curves to illustrate ac-
     curacy improvement with increasing sample
     size.

Fig. 5 shows our proposed methodologies flowchart. In
the end, by demonstrating the efficacy of ML models and FS
methods on the Bonn dataset, this study advances the area
of epilepsy prediction. The findings highlight the potential of
these approaches in developing reliable and accurate prediction
systems for epilepsy management.



Fig. 5. Overall architecture of proposed methodology.

### C. Setup

The suggested system requires importing numpy, pandas,
scikit-learn, matplotlib, xgboost, and lightgbm. The target
variable is preprocessed after loading the CSV dataset. Detect
and eliminate outliers, then scale features using "Standard-
Scaler". Features are chosen using RFE. After separating the
dataset into training and test sets using the 'train_test_split'
function, we trained our models using the set. Next, we
assessed accuracy, sensitivity, specificity, precision, recall, and
F1-score. Then performance comparison has done to show the
practicality of the proposed hybrid model.

### D. EpiNet

Within the suggested approach, a simple averaging tech-
nique is used to build the hybrid model, named EpiNet.
After training multiple machine learning models, including
SVM, Gaussian Naive Bayes, Gradient Boosting, XGBoost,
and LightGBM, the average predictions are pooled. This
cooperative strategy uses each model's advantages to ensure
fair and efficient decision-making. For each instance, the
hybrid model's output is the average forecast from the group's
insights, improving epileptic seizure prediction accuracy and
reliability. The hybrid model is built using simple averaging
in the provided way. After training multiple machine learning
models, including SVM, Gaussian Naive Bayes, Gradient
Boosting, XGBoost, and LightGBM, the average predictions
are pooled. This cooperative strategy uses each model's ad-
vantages to ensure fair and efficient decision-making. The
Hybrid model's final output for a given instance is the average
forecast from the group's insights, improving epileptic seizure
prediction accuracy and reliability.

## VI. RESULTS AND DISCUSSION

Our extensive epileptic episode prediction study comprised
machine learning models and innovative feature selection
methods. This method validated each model and explained
seizure prediction. This discussion will explain the findings
and any surprises. Linking our findings to earlier research aids
comprehension. Comprehensive analysis is needed to improve
epileptic seizure prediction.

### A. Model Performance Dissection

To evaluate models' prediction ability, the suggested tech-
nique uses major performance metrics. Model accuracy, sensi-
tivity, specificity, precision, recall, and F1-score show perfor-
mance. The accuracy metric the ratio of accurate predictions
to total forecasts indicates correctness. Recall, or sensitivity,

measures how effectively models detect positive cases. Specificity, which assesses negative projections, boosts sensitivity. Recall retrieves all positive events, whereas precision assesses positive prediction accuracy. The balanced F1-score combines recollection and accuracy for a detailed evaluation. The mathematical formulations of these metrics' equations give a solid foundation for assessing model performance in varied circumstances. The formulaś used to calculate the evaluation metrics are as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{1}$$

$$Sensitivity(Recall) = \frac{TP}{TP + FN} \tag{2}$$

$$Specificity = \frac{TN}{TN + FP} \tag{3}$$

$$Precision = \frac{TP}{TP + FP} \tag{4}$$

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall} = \frac{2 * TP}{2 * TP + FP + FN} \tag{5}$$

TABLE II. PERFORMANCE ANALYSIS OF THE PROPOSED EPINET MODEL

| Model | Sensitivity | Specificity | Precision | Recall | F1-Score | Support | ROC AUC |
|---|---|---|---|---|---|---|---|
| SVM | 0.909 | 1.000 | 0.998 | 0.998 | 0.998 | 22 | 0.995 |
| GaussianNB | 0.954 | 0.983 | 0.989 | 0.998 | 0.984 | 22 | 0.993 |
| Gradient Boosting Classifier | 0.909 | 0.999 | 0.997 | 0.997 | 0.997 | 22 | 0.974 |
| XGB Classifier | 0.909 | 1.000 | 0.998 | 0.998 | 0.998 | 22 | 0.996 |
| LGBM Classifier | 0.909 | 1.000 | 0.998 | 0.998 | 0.998 | 22 | 0.985 |
| **Hybrid** | **0.909** | **1.000** | **0.998** | **0.998** | **0.998** | **22** | **0.971** |

Our method is evaluated using SVM, Gaussian Naive Bayes, Gradient Boosting, XGBoost, LightGBM, and the innovative hybrid model in Table II. Each epileptic seizure prediction method has merits and downsides. SVM and XGBoost excel in accuracy and sensitivity. They recognise non-seizure circumstances well. For seizure prediction, Gaussian Naive Bayes and Gradient Boosting are more sensitive. Seizure treatment relies on their sensitivity to detect true positives. We now notice the hybrid model (EpiNet) for another reason. It balances several variables well. Its strength is combining the best features of various models without compromising others. It predicts epileptic seizures well due to its comprehensive approach. This hybrid model's success raises questions regarding its components' interactions. Exploring each model's strengths reveals their goal: accurate seizure prediction. Details on the hybrid model reveal its effectiveness and probable linkages, raising interest in its role in epilepsy forecasting.

## B. Insights into Hybrid Model Superiority

Table III shows the mixed model ratings where the evaluation metrices such as accuracy, precision, recall, and F1-score of Class 0 (Non-Epilepsy) and 1 (Epilepsy) are measured. Also macro and weighted averages are provided. The hybrid model predicts epileptic episodes with 99.81% accuracy. The Hybrid Model adeptly amalgamates predictions from diverse models, showcasing superior performance in specificity and precision compared to individual models. Simultaneously, it achieves heightened sensitivity, signifying its efficacy in achieving an optimal balance. This amalgamation strategically leverages the strengths of individual models, providing a robust and superior framework for epileptic seizure prediction. Importantly, our findings demonstrate that the Hybrid Model outperforms individual models in key metrics, aligning with emerging research [3] advocating for hybridization to substantially enhance predictive accuracy. Our findings echo and extend the conclusions drawn by prior studies [21] that emphasized the significance of feature selection techniques in enhancing predictive accuracy. The nuanced performance variations observed in SVM models align with previous assertions regarding the impact of dataset characteristics on SVM effectiveness [22]. Additionally, the superior performance of the hybrid model resonates with recent research advocating for ensemble methods [10] in achieving superior predictive accuracy.

TABLE III. EVALUATION METRICES FOR THE EPINET MODEL

| Metric | Precision | Recall | F1-Score | Accuracy |
|---|---|---|---|---|
| Class 0 (Non-Epilepsy) | 1.00 | 1.00 | 1.00 | **99.81%** |
| Class 1 (Epilepsy) | 1.00 | 0.91 | 0.95 | |
| **Macro Avg** | 1.00 | 0.95 | 0.98 | **99.81%** |
| **Weighted Avg** | 1.00 | 1.00 | 1.00 | |



Fig. 6. Accuracy curve.

As shown in Fig. 6, the relationship among the amount of specimens and the precision of the predictions is represented by the accuracy curve. It illustrates how the precision of the models increases with the inclusion of additional samples. Sensitivity, denoted by the ROC AUC curve in (see Fig. 7) is the compromise between the FPR and the TPR. The graphical depiction illustrates the capacity of the models to distinguish positive from negative classifications. An increased AUC (Area Under the Curve) signifies superior data classifi-

Fig. 7. ROC AUC curve.



Fig. 8. Confusion matrix.

cation performance. The ROC curve analysis shows not only the exceptional classification performance of individual models and EpiNet but also indicates their specificity and sensitivity balance. The perfect AUC score of 1.00 for EpiNet suggests no overlap between the true positive rate and false positive rate, indicating an ideal separation of classes. The hybrid model's confusion matrix (see Fig. 8) shows that the vertical elements show the true positive and true negative forecasts, which show cases that were correctly labelled as Non-Epilepsy and Epilepsy, respectively. Specifically, the model's perfect success rate of 100% for non-epilepsy instances demonstrated its exceptional ability to identify such situations. For the Epilepsy class, the model's sensitivity was 90.91%, indicating that it accurately detected 90.91% of actual instances. In contrast, the model's ability to accurately identify genuine negative instances is shown by the 100% accuracy for the Non-Epilepsy class. On the accuracy curve, EpiNet maintains a plateau of high accuracy across sample sizes, demonstrating robustness against overfitting—a challenge often encountered with smaller datasets. The nuanced fluctuations observed in other models' accuracy curves could be attributed to their individual handling of the dataset's complexity, which is mitigated in EpiNet's ensemble strategy. These insights into the model performance dynamics affirm the superiority of the hybrid approach, where collective intelligence effectively captures the intricate patterns in EEG data critical for seizure prognosis.

The hybrid model distinguishes epilepsy from non-epilepsy. This prepares for epilepsy prediction and classification. Our method diagnoses epilepsy patients with 99.81% accuracy. The findings demonstrate that our epilepsy prediction method is accurate and reliable. This confirms our model's applicability.

From Table IV, we can see that the EpiNet Model, a combination of different techniques for epileptic seizure prediction, outperforms individual Support Vector Machines (SVM) models, achieving an accuracy of 99.81%. Given the scarcity of studies on hybrid models, our approach, incorporating SVM alongside other methods, demonstrates enhanced prediction accuracy. Overall, our system, leveraging various machine learning models, presents promising results, with the Hybrid

TABLE IV. PERFORMANCE ANALYSIS OF THE PROPOSED EPiNET MODEL WITH STATE-OF-ARTS BASED ON SVM MODEL

| Publication | Models | Accuracy |
|---|---|---|
| [26] | SVM | 94% |
| [21] | "SVM" with "LBP" | 97.80% |
| [10] | SVM | 92.23% |
| [30] | SVM | 97% |
| [23] | SVM | 95.33% |
| [22] | SVM | 94% |
| **Proposed Method** | **Hybrid Model** | **99.81%** |

Model showing superior performance and accuracy compared to individual models.

### C. Identification of Model-Specific Strengths

Delving deeper, the unexpected prominence of specific models in certain metrics prompts a nuanced understanding. SVM emerges as a stalwart in specificity, a characteristic well-documented in studies focusing on Support Vector Machines for epilepsy detection [26]. XGBoost, with its gradient boosting prowess, excels in overall accuracy, an attribute substantiated in recent research [38]. Identifying these model-specific strengths contributes valuable insights for tailored model selection based on the diagnostic emphasis.

### VII. CONCLUSION AND FUTURE WORK

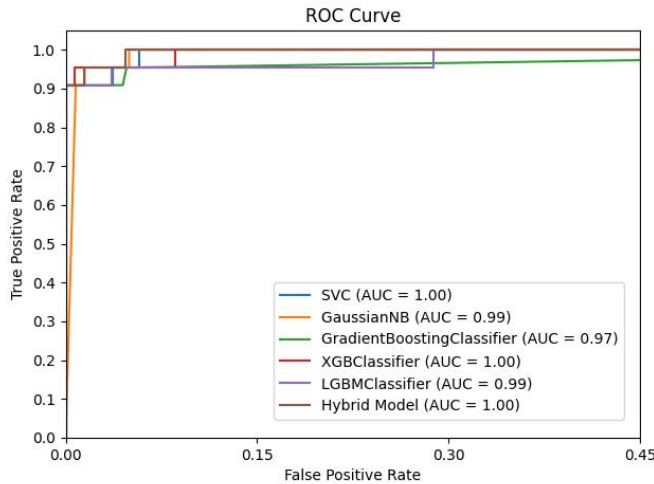In the pursuit of advancing automatic epilepsy detection from EEG signal data, this study introduces a novel hybrid machine learning model, exhibiting a remarkable accuracy of 99.81% in the classification of five distinct classes (A-B-C-D-E) based on a dataset from the University of Bonn. In this study, we carefully prepared EEG data, addressed outliers, scaled features, and trained our model. What sets our approach apart is the use of a straightforward averaging method to combine model predictions, yielding outstanding classification results. This work establishes a strong foundation for advancing accurate epileptic seizure prediction techniques. By contributing to the landscape of medical diagnostics, this research holds the potential to significantly enhance healthcare strategies tailored to epilepsy patients.

While this study demonstrates promising results in automatic epilepsy detection, it is essential to acknowledge the limitations stemming from dataset specificity and potential challenges in generalization. Though our model excelled on the Bonn dataset, its generalizability to diverse datasets requires careful consideration. It signifies a need for further exploration and adaptation to diverse data sources. Adapting and fine-tuning the model for consistent performance across different data sources is necessary. Our findings in epileptic seizure prediction using EpiNet open new avenues in proactive healthcare management. By achieving a high accuracy rate, this model can be integrated into real-time monitoring systems, offering early warning signals for impending seizures. This advancement holds the potential to drastically improve patient outcomes, enabling timely interventions and personalized treatment plans. Future work will focus on enhancing the generalizability of EpiNet across diverse datasets, ensuring its applicability in various clinical settings. By doing so, we aim to contribute significantly to the evolution of healthcare strategies for epilepsy patients, making a tangible difference in their quality of life and care. Future research endeavors should focus on addressing these constraints by exploring diverse datasets, refining model robustness, and incorporating real-time monitoring systems. Overcoming these challenges will contribute to the continuous evolution of accurate and reliable epileptic seizure prediction techniques, further advancing the field of medical diagnostics and improving healthcare for individuals with epilepsy.

### References

[1] https://www.who.int/news-room/fact-sheets/detail/epilepsy (Accessed on: 17 January 2024).

[2] O. K. Esha, N. Haque, F. Ahmed and N. Begum, "A Comprehensive Review on Epilepsy Diagnosis and Prognosis using Machine Learning and Deep Learning Approaches," 2022 IEEE International Women in Engineering (WIE) Conference on Electrical and Computer Engineering (WIECON-ECE), Naya Raipur, India, 2022, pp. 159-164, doi: 10.1109/WIECON-ECE57977.2022.10151006.

[3] Kapoor, Bhaskar, et al. "Epileptic Seizure Prediction Based on Hybrid Seek Optimization Tuned Ensemble Classifier Using EEG Signals." Sensors, vol. 23, no. 1, Dec. 2022, p. 423. Crossref, https://doi.org/10.3390/s23010423.

[4] Savadkoohi, Marzieh, Timothy Oladunni, and Lara Thompson. "A machine learning approach to epileptic seizure prediction using Electroencephalogram (EEG) Signal." Biocybernetics and Biomedical Engineering 40.3 (2020): 1328-1341.

[5] Slimen, Itaf Ben, Larbi Boubchir, and Hassene Seddik. "Epileptic seizure prediction based on EEG spikes detection of ictal-preictal states." Journal of biomedical research 34.3 (2020): 162.

[6] Usman, Syed Muhammad, Shehzad Khalid, and Zafar Bashir. "Epileptic seizure prediction using scalp electroencephalogram signals." Biocybernetics and Biomedical Engineering 41.1 (2021): 211-220.

[7] Usman, Syed Muhammad, Shehzad Khalid, and Muhammad Haseeb Aslam. "Epileptic seizures prediction using deep learning techniques." Ieee Access 8 (2020): 39998-40007.

[8] Alshebeili, Saleh A., et al. "Inspection of EEG signals for efficient seizure prediction." Applied Acoustics 166 (2020): 107327

[9] Wang, Gang, et al. "Seizure prediction using directed transfer function and convolution neural network on intracranial EEG." IEEE Transactions on Neural Systems and Rehabilitation Engineering 28.12 (2020): 2711-2720.

[10] Ryu S, Joe I. A Hybrid DenseNet-LSTM Model for Epileptic Seizure Prediction. Applied Sciences. 2021; 11(16):7661. https://doi.org/10.3390/app11167661https://doi.org/10.3390/app11167661.

[11] Viana, Pedro F., et al. "Seizure forecasting using minimally invasive, ultra-long-term subcutaneous electroencephalography: Individualized intrapatient models." Epilepsia (2022).

[12] Behnoush, B., et al. "Machine learning algorithms to predict seizure due to acute tramadol poisoning." Human & Experimental Toxicology 40.8 (2021): 1225-1233.

[13] Syed Muhammad Usman, Muhammad Usman, Simon Fong, "Epileptic Seizures Prediction Using Machine Learning Methods", Computational and Mathematical Methods in Medicine, vol. 2017, Article ID 9074759, 10 pages, 2017.https://doi.org/10.1155/2017/9074759.

[14] Echtioui A, Hamam H,et al. Deep Learning Models for Predicting Epileptic Seizures Using iEEG Signals. Electronics. 2022; 11(4):605. https://doi.org/10.3390/electronics11040605

[15] Harpale, Varsha, and Vinayak Bairagi. "An adaptive method for feature selection and extraction for classification of epileptic EEG signal in significant states." Journal of King Saud University-Computer and Information Sciences 33.6 (2021): 668-676.

[16] Tran, Ly V., et al. "Application of machine learning in epileptic seizure detection." Diagnostics 12.11 (2022): 2879

[17] Cao, Xincheng, et al. "Automatic seizure classification based on domain-invariant deep representation of EEG." Frontiers in Neuroscience 15 (2021): 760987.

[18] Tripathi, Prashant Mani, et al. "Automatic Seizure Detection and Classification Using Super-resolution Superlet Transform and Deep Neural Network-A Preprocessing-less Method." Computer Methods and Programs in Biomedicine (2023): 107680.

[19] Ahmed, Mohammed Imran Basheer et al. "A Review on Machine Learning Approaches in Identification of Pediatric Epilepsy." SN computer science vol. 3,6 (2022): 437. doi:10.1007/s42979-022-01358-9

[20] Farooq, Muhammad Shoaib, Aimen Zulfiqar, and Shamyla Riaz. "Epileptic Seizure Detection Using Machine Learning: Taxonomy, Opportunities, and Challenges." Diagnostics 13.6 (2023): 1058.

[21] Tharayil JJ, Chiang S, Moss R, Stern JM, Theodore WH, Goldenholz DM. A big data approach to the development of mixed-effects models for seizure count data. Epilepsia. 2017;58(5):835-844. doi:10.1111/epi.13727

[22] Jaiswal, A.K., & Banka, H. (2017). Epileptic seizure detection in EEG signal with GModPCA and support vector machine. Bio-medical materials and engineering, 28 2, 141-157 .

[23] Tiwari, A.K., Pachori, R.B., Kanhangad, V., & Panigrahi, B.K. (2017). Automated Diagnosis of Epilepsy Using Key-Point-Based Local Binary Pattern of EEG Signals. IEEE Journal of Biomedical and Health Informatics, 21, 888-896.

[24] Toraman, Suat. "Preictal and Interictal Recognition for Epileptic Seizure Prediction Using Pre-trained 2DCNN Models." Traitement du Signal 37.6 (2020).

[25] Satapathy, Sandeep & Jagadev, Alok & Dehuri, Satchidananda. (2017). Weighted Majority Voting Based Ensemble of Classifiers Using Different Machine Learning Techniques for Classification of EEG Signal to Detect Epileptic Seizure. Informatica. 41. 99-110.

[26] Lima, C.A., Coelho, A.L., Madeo, R.C., & Peres, S.M. (2015). Classification of electromyography signals using relevance vector machines and fractal dimension. Neural Computing and Applications, 27, 791 - 804.

[27] Kabir, E., Siuly, & Zhang, Y. (2016). Epileptic seizure detection from EEG signals using logistic model trees. Brain informatics, 3(2), 93–100. https://doi.org/10.1007/s40708-015-0030-2

[28] Aslam, M.H.; Usman, S.M.; Khalid, S.; Anwar, A.; Alroobaea, R.; Hussain, S.; Almotiri, J.; Ullah, S.S.; Yasin, A. Classification of EEG Signals for Prediction of Epileptic Seizures. Appl. Sci. 2022, 12, 7251. https://doi.org/10.3390/ app12147251

[29] Patrick H, Luckett BS. "Nonlinear methods for detection and prediction of epileptic seizures." A Dissertation submitted in University of South Alabama, July, 2018.

[30] Abualsaud, Khalid et al. "Ensemble classifier for epileptic seizure detection for imperfect EEG data." TheScientificWorldJournal vol. 2015 (2015): 945689. doi:10.1155/2015/945689

[31] Sharma, Rahul. "Localization of epileptic surgical area using automated hybrid approach based on higher-order statistics with sensitivity analysis and residual wavelet transform." Biomedical Signal Processing and Control 86 (2023): 105192.

[32] Brari, Zayneb, and Safya Belghith. 'A Novel Machine Learning Approach for Epilepsy Diagnosis Using EEG Signals Based on Correlation Dimension'. IFAC-PapersOnLine, vol. 54, no. 17, 2021, pp. 7–11, https://doi.org10.1016/j.ifacol.2021.11.018.

[33] Slimen IB, Boubchir L, Seddik H. Epileptic seizure prediction based on EEG spikes detection of ictal-preictal states. J Biomed Res. 2020 Feb 17;34(3):162-169. doi: 10.7555/JBR.34.20190097. PMID: 32561696; PMCID: PMC7324272.

[34] Ilakiyaselvan, N et al. "Deep learning approach to detect seizure using reconstructed phase space images." Journal of biomedical research vol. 34,3 (2020): 240-250. doi:10.7555/JBR.34.20190043

[35] Singh K, Malhotra J. Two-layer LSTM network-based prediction of epileptic seizures using EEG spectral features. Complex Intell Syst. 2022;8(3):2405-2418. doi:10.1007/s40747-021-00627-z

[36] Li, Chang, et al. 'EEG-Based Seizure Prediction via Transformer Guided CNN'. Measurement, vol. 203, 2022, p. 111948, https://doi.org10.1016/j.measurement.2022.111948.

[37] Rasheed, Khansa et al. "A Generative Model to Synthesize EEG Data for Epileptic Seizure Prediction." IEEE transactions on neural systems and rehabilitation engineering : a publication of the IEEE Engineering in Medicine and Biology Society vol. 29 (2021): 2322-2332. doi:10.1109/TNSRE.2021.3125023

[38] Kumar, T., et al. 'A Modified XG Boost Classifier Model for Detection of Seizures and Non-Seizures'. WSEAS TRANSACTIONS ON BIOLOGY AND BIOMEDICINE, vol. 19, 01 2022, pp. 14–21, https://doi.org10.37394/23208.2022.19.3.

[39] Usman, Syed Muhammad, Shehzad Khalid, and Sadaf Bashir. "A deep learning based ensemble learning method for epileptic seizure prediction." Computers in Biology and Medicine 136 (2021): 104710.

[40] "Bonn Dataset" https://www.ukbonn.de/epileptologie/arbeitsgruppen/ag-lehnertz-neurophysik/downloads/

# A Comparative Study of ChatGPT-based and Hybrid Parser-based Sentence Parsing Methods for Semantic Graph-based Induction

Walelign Tewabe[1], László Kovács[2]

Institute of Technology, University of Miskolc, Hungary[1,2]

Institute of Technology, Debre Markos University, Ethiopia[1]

*Abstract*—Sentence parsing is a fundamental step in the conversion of a text document into semantic graphs. In this research, novel phrase parsing techniques for semantic graph-based induction are presented, namely the ChatGPT-based and Hybrid Parser-based approaches. The performance of these two approaches in the context of inducing semantic networks from textual data is assessed through a comprehensive analysis in this study. The primary purpose is to enhance the construction of semantic graphs, specifically focusing on capturing detailed event descriptions and relationships within text. The research finds that the Hybrid Parser-Based approach exhibits a slight advantage in accuracy (acc_hybrid = 0.87) compared to ChatGPT (acc_GPT = 0.85) in sentence parsing tasks. Furthermore, the efficiency analysis reveals that ChatGPT's response quality varies with different prompt sizes, while the Hybrid Parser-Based method consistently maintains an "excellent" response quality rating.

*Keywords*—*Adverb prediction; ChatGPT; hybrid parser-based; natural language processing; sentence parsing; semantic graph induction*

## I. INTRODUCTION

Semantic Graph Induction is a computational approach in Natural Language Processing (NLP) and artificial intelligence that aims to extract and represent structured knowledge and semantic relationships from unstructured textual data. Semantic Graph visually represents the semantic structure of a document extracted from sentences [1]. Semantic graphs play a multifaceted role in various applications, spanning information retrieval, knowledge representation, question answering, text summarization, document clustering, and classification.

Furthermore, in the finance industry, semantic graphs have emerged as a crucial tool for managing financial knowledge securely [11], enabling applications like transaction surveillance, financial crime detection and prevention, and non-compliant user detection [12]. In the entertainment industry, particularly social media, knowledge graphs power social graphs that help platforms like Facebook connect users within the context of their relationships, while also enhancing recommender systems to offer personalized content recommendations based on user interests [13]. Moreover, semantic graphs play a vital role in cybersecurity by mapping historical cyber attacks and predicting potential future breaches, thus bolstering cyber defense strategies [14].

This study explores the creation of semantic graphs, which are visual representations of knowledge and the interconnections between concepts. Specific tools within the domain of NLP parsing are working for constructing these semantic graphs. However, there are limitations in their ability to present detailed event descriptions, particularly concerning time and place. Recognizing the limitations present in current NLP parsing tools, the primary objective of this research is to enhance the existing approach. To address these limitations, this paper introduces a solution that involves identifying all functional components, including Subject, Predicate, Direct Object, Indirect Object, and Conjunction. Simultaneously, the method explores the prediction of adverb types, encompassing Time, Place, Manner, Degree, and Frequency, thus enriching the depth of linguistic analysis.

To gain a deeper understanding of knowledge, concepts, and the complex web of relationships between them, this research extends beyond traditional limitations by incorporating a more comprehensive set of components. Specifically, the study introduces novel ChatGPT-based and Hybrid Parser-based Semantic Graph Construction and conducts a comparative analysis. This analysis assesses the details of these two approaches, dissecting their respective strengths, weaknesses, and applications.

In this regard, ChatGPT is one of the state-of-the-art Large Language Models (LLMs) [15], that has emerged as a transformative force in the field of NLP. It plays a pivotal role in the construction of semantic graphs by leveraging their natural language understanding capabilities. These models are trained on extensive text corpora and can extract and encode intricate relationships between concepts and entities within textual data. ChatGPT's previous experiences with these tasks are informed by its extensive pre-training on a diverse range of internet text [16]. This pre-training allows it to understand and generate human-like text and perform tasks related to semantic graph construction with high accuracy. By leveraging this understanding, ChatGPT can contribute significantly to the creation and enrichment of semantic graphs across various domains, from healthcare [10] and finance to information retrieval and content recommendation [17]. It has demonstrated remarkable skill in a wide array of language understanding tasks, including question-answering, language generation, and text summarization [18]. However, the question arises: can ChatGPT be effectively harnessed to tackle the difficulties of semantic graph-based induction? On the other hand, Hybrid Parser-based methods integrate multiple NLP components, combining rule-based and machine-learning techniques, to extract and represent semantic relationships from text. The marriage of these disparate approaches promises enhanced

robustness and adaptability. This study sets out to investigate which of these approaches outshines in the domain of semantic graph construction, and whether a hybrid approach provides a balanced solution. The contributions of this work are Semantic Graph Construction Enhancement, Testing a novel ChatGPT-based parsing for functional sentence parsing, and Comparative Analysis of Methodologies.

The paper is structured as follows: In the second section, a semantic graph construction model is presented, and a detailed procedure for building the presented model is provided. We discuss the latest NLP background technology and results. Additionally, we explore different knowledge base resources and their applications. The third section describes the proposed Hybrid Parser-based method, explaining all process steps. In the fourth section, we describe the ChatGPT-based method, encompassing the environment, dataset size, benchmark, and evaluation methods. Next, we present the experimental results, analyze the evaluation findings from multiple perspectives, and demonstrate the potential applications of our approach. In the sixth section, we conduct an efficiency analysis and engage in a discussion. Finally, in the concluding section, we summarize our findings and offer suggestions for future research directions.

## II. LITERATURE REVIEW

### A. Basics of Semantic Graph

A semantic graph is a graph model where nodes represent concepts and edges (or arcs) represent relationships between those concepts [19]. This model type is often used in artificial intelligence applications for representing knowledge.

### B. Definition 2.1

A graph $G = (V, E)$ is defined by a set of nodes $V$ and a set of edges $E$ between these nodes. Let $E \subseteq V \times V$ represent directed edges or arcs [20]. Each directed edge $(u, v) \in E$ signifies a connection from a start (tail) vertex $u$ to an end (head) vertex $v$, where $u$ and $v$ are elements of the node set $V$. The graph's structure is characterized by these directed connections, providing a representation of relationships between nodes. Each node is associated with a label $Label(v)$.

Building semantic graphs is essential for many practical uses and ongoing research [8], [21], [22]. As we have more and more data available, creating these meaningful graphs becomes increasingly important for learning from different sources. Scientists keep looking for new ways to make this field better, and they use it in things like understanding language, organizing knowledge, and using artificial intelligence. They make structured graphs and networks to show how words, ideas, and things are connected. These graphs help in finding information, answering questions, and suggesting things you might like. So, making these graphs is a big part of helping computers and people work together better. When texts are represented graphically, it allows the preservation of additional information like the text's inner structures, semantic relationships, and term order. However, events like these are not effectively captured using current NLP parsing and semantic graph construction. As an illustration, Fig. 1 provides a visual



Fig. 1. A visualization of the basic event knowledge graph for eating [9].

insight into a fundamental event knowledge graph centered around the concept of "eating" [9].

Understanding natural language is a big challenge, and that's where semantic graphs come into play. Enhancing our grasp of natural language relies heavily on the development of semantic graphs, a field that's been increasingly in the spotlight. Researchers are actively exploring the creation of these graphs and how they can represent knowledge, diving into structured data, relationships, and more detailed elements, which align with prior work on Semantic Role Labeling (SRL) and adverb sense disambiguation. These efforts aim to provide a more comprehensive understanding of semantic parsing, event descriptions, and the complexities involved, as outlined in related works [23].

Knowledge graphs have also got substantial attention in recent years, serving as vital tools for organizing and connecting vast amounts of information from diverse sources, including text corpora, databases, and the web [24]. Some well-known knowledge graphs, such as DBpedia, Freebase, and Wikidata, have been crucial in this effort. We're also using some smart techniques like word embeddings and word vector representations to make semantic graphs even better [25].

Resource Description Framework (RDF) and ontologies are the foundation for constructing structured, machine-readable semantic graphs, playing a pivotal role in knowledge representation and the advancement of the semantic web. RDF, with its subject-predicate-object triples and Uniform Resource Identifiers (URIs), ensures global consistency and interoperability. Ontologies, including OWL and RDFS, enrich RDF's capabilities by defining the vocabulary and structure for resources and relationships within specific domains, making it easier to understand and work with the information [36]. Together, RDF and ontologies are super important for making and using semantic graphs across different fields.

At the same time, the Semantic Web initiative is pushing for structured data to be shared and linked on the web. They're

using things like Linked Data, RDF, and SPARQL Protocol and RDF Query Language (SPARQL) queries to create big semantic graphs that cover a lot of the web[37]. But there are challenges too. We need better ways to handle big sets of data, put together text and visual data, and make sure the knowledge graphs we create are complete and correct. Researchers are used new techniques, like word embeddings and entity embeddings, to help to understand the fine details of how words and things are related [7]. As we have more and more data, making meaningful semantic graphs becomes super important for getting useful information from different places.

In general, the fields of RDF, ontologies, and the ideas behind the Semantic Web initiative where semantic graph play an important role that understand and manage information. The semantic graphs serve as a crucial foundation for knowledge representation and data integration, facilitating the consistence management of structured data on the web. However, this field is evolving, with ongoing efforts focused on improving graph construction techniques, addressing data handling challenges, and harnessing the power of embedding techniques to capture richer semantic relationships. As the landscape of available data continues to expand, the construction of semantic graphs becomes essential for unlocking valuable insights and enabling data-driven applications across various domains.

### III. ChatGPT-based Sentence Parsing

A significant aspect of language models is the LLM, recognized for its capacity to achieve a wide-ranging understanding of language and proficiently generate text. LLMs acquire this capability through an extensive training process where they learn from vast amounts of data, effectively processing billions of parameters. This training demands substantial computational resources [28]. These language models primarily employ artificial neural networks, predominantly relying on transformer architectures, and undergo (pre-)training utilizing self-supervised and semi-supervised learning approaches [29].

Functioning as autoregressive language models, LLMs operate by taking an input text and iteratively predicting subsequent tokens or words [30]. Until the year 2020, the primary approach to adapt these models for specific tasks was fine-tuning. However, with the emergence of larger models like GPT-3, they can now be engineered with prompts to achieve similar outcomes [31]. LLMs are believed to acquire an inherent understanding of syntax, semantics, and the "ontology" within human language corpora [32].

Prominent examples of LLMs include OpenAI's GPT models like Generative Pre-trained Transformer (GPT)-3.5 and GPT-4, Google's Pathways Language Model (PaLM) employed in Bard, Meta's Language Model for Language Modeling (LLaMa), as well as BigScience Large Open-science Open-access Multilingual Language Model (BLOOM), Ernie 3.0 Titan, and Anthropic's Claude 2. In this study, due to the model's capabilities, researchers utilized the ChatGPT 3.5 OpenAI API for the sentence parsing.

#### A. Basics of GPT-based Models

Chat GPT (Generative Pre-trained Transformer) models are designed to understand and generate human-like text by processing vast amounts of data during training [34]. They operate by predicting the next word in a sequence of words and have been instrumental in various NLP tasks. Understanding these fundamental concepts is essential for harnessing the power of GPT-based models in language-related applications. The accuracy of the ChatGPT 3.5 model heavily relies on the quality and representativeness of the labeled dataset used for fine-tuning [35]. The pre-trained ChatGPT model is fine-tuned on a labeled dataset of adverbs to improve its categorization accuracy.

#### B. The Architecture of ChatGPT

ChatGPT is based on the transformer architecture, that allows for parallel processing, which makes it well-suited for processing sequences of data such as text. ChatGPT uses the PyTorch library, an open-source machine learning library, for implementation. ChatGPT is made up of a series of layers, each of which performs a specific task.

#### C. Prompt Engineering Techniques

Prompt engineering is a crucial technique employed to guide the behavior of large-scale language models like ChatGPT [34]. By strategically constructing input prompts, researchers and developers aim to obtain more accurate and relevant responses from these models [33]. Several prompts engineering strategies, including prompt rewriting, contextual incorporation, explicit instructions, and templates, have been proposed to address control and responsiveness challenges, aligning the model's outputs with user targets and expectations. The careful design of prompts plays a pivotal role in influencing the quality and relevance of ChatGPT's responses, making it a valuable skill for those working with AI systems. For instance, in a real-world context, prompt engineering bears the potential to enhance the efficiency, accuracy, and effectiveness of healthcare delivery by guiding AI models to provide valuable insights and solutions. However, it's crucial to acknowledge the limitations and risks associated with AI, such as the model's inability to access real-time data or offer personalized medical advice. This necessitates verification by qualified professionals and raises concerns about privacy and data security. Despite these challenges, the significance of prompt engineering has seen exponential growth since the inception of ChatGPT, with ongoing research endeavors aimed at refining and expanding this critical skill, particularly within the medical field. In this specific study, researchers have developed and employed high-quality training sets as templates for prompts to augment the accuracy of responses.

#### D. Methodology

The methodology for this study involves the following steps.

*1) Construction of a Labeled Dataset:* A high-quality labeled dataset is carefully collected to fine-tune ChatGPT for sentence parsing by including the adverb type prediction. This dataset includes Subject, Predicate, Direct Object, Indirect Object, Conjunction, and adverb types such as Time, Place, Manner, Degree, and Frequency. The dataset is essential for training ChatGPT to categorize adverbs accurately and for sentence parsing.

```
# Prompt for generating a semantic graph description
prompt_template= """
Sentence 1: The coffee shop is always busy in the morning.
Parsing Answer 1: {'Predicate': is, 'Subject': The coffee shop,
    'Direct Object': [], 'Indirect Object': [], 'Time': in the
    morning, 'Place': [], 'Manner': always busy, 'Frequency': [],
    'Degree': []}
Sentence 2: The train arrived at the station on time.
    Parsing Answer 2: {'Predicate': arrived, 'Subject': The train,
    'Direct Object': [], 'Indirect Object': [], 'Time': on time,
    'Place': at the station, 'Manner': [], 'Frequency': [], 'Degree':
    []}
Sentence 3: Ethiopia defeated Italy at the Battle of Adwa.
    Parsing Answer 2: {'Predicate': defeated, 'Subject': Ethiopia,
    'Direct Object': Italy, 'Indirect Object': [], 'Time': [],
    'Place': at the Battle of Adwa, 'Manner': [], 'Frequency': [],
    'Degree': []} """

custom_prompt=""" Generate the Predicate, Subject, Direct Object,
Indirect Object, Time, Place, Manner, Frequency, Degree, Conjunction,
Clause parts of the following sentence:- Sentence: 'The child reads
the book carefully and attentively  at the library everyday. """

prompt = prompt_template + custom_prompt
```

Fig. 2. Sample prompt template [9].

*2) Fine-tuning ChatGPT:* Fine-tuning is a phase where the pre-trained model is further trained on the specific task it will be used for. The objective of this phase is to adapt the model to the specific task and fine-tune the parameters so that the model can produce outputs that are in line with the expected results. The pre-trained ChatGPT 3.5 model is fine-tuned using the labeled dataset of functional sentence structure. One of the most important things in the fine-tuning phase is the selection of the appropriate prompts. The prompt is the text given to the model to start generating the output. Providing the correct prompt is essential because it sets the context for the model and guides it to generate the expected output. It is also important to use the appropriate parameters during fine-tuning, such as the temperature, which affects the unpredictability of the output generated by the model. As shown in Fig. 2 the researcher developed and used representative prompt templates from the collected dataset in this regard. This fine-tuning process helps the model to learn and recognize the functional structure of a sentence including the adverb types based on contextual information.

*3) Response Generation:* With the ability to predict the functional structure of the sentence, ChatGPT can generate coherent and contextually relevant responses. These responses are informed by the adverb-type predictions, making them more precise and contextually appropriate. Throughout the methodology, emphasis is placed on the quality and representativeness of the labeled dataset, as this significantly influences the accuracy of adverb categorization and response generation. This ChatGPT-based methodology combines the power of pre-trained language models with fine-tuning on a domain-specific dataset to enhance adverb type prediction and response generation. It is a dynamic approach that leverages ChatGPT's natural language understanding and generation capabilities, making it a valuable tool for various NLP applications.

## IV. HYBRID PARSER-BASED METHOD

The creation of a Hybrid Parser-based sentence parsing framework is a noteworthy breakthrough in the field of NLP.



Fig. 3. The structural framework of the proposed hybrid parser based sentence parsing.

This innovative approach combines rule-based and machine-learning methods to extract meaning from text [38], addressing the limitations of current NLP parsing techniques. By incorporating both rule-based and machine-learning components, this framework becomes capable of handling a wider range of linguistic structures and domains, ensuring robust performance. Its primary objective is to enhance the accuracy of semantic parsing by capturing context-specific elements in language, ultimately improving the comprehension of the underlying meaning in the text. The framework strikes a careful balance between accuracy and efficiency, allowing for the precise construction of a semantic graph from textual content. The architecture of this framework encompasses text preprocessing, rule-based and machine learning-based sentence parsing, adverb-type prediction, and semantic graph construction.

One distinguishing feature of this framework is its dedicated component for predicting adverb types within the text. This feature plays a pivotal role in accurately extracting the essence of a sentence. The integration of outputs from both rule-based and machine learning-based parsing yields a comprehensive semantic graph representing the structured knowledge present in the text. This Hybrid parser-based approach harnesses the strengths of rule-based systems, which excel at handling linguistic patterns and prior knowledge, and machine learning models, which adapt to context and data-driven insights. As a result, the framework enhances natural language understanding and information extraction, offering a promising solution to the challenges presented by traditional parsing methods.

Fig. 3 provides an overview of the structural framework of the Hybrid Parser-based approach. It illustrates the key components, including text preprocessing, rule-based and machine learning-based parsing, adverb-type prediction, and semantic graph construction, highlighting their interconnections.

### A. Methodology

The researchers utilized a free cloud-based platform called Google Collaboratory for running and writing Python code. For text analysis and parsing, we used essential parsing tools such as spaCy and NLTK. To improve the analysis and understanding of language, we integrated external resources, including dictionaries like Webster and ontologies such as WordNet.

Furthermore, to train the adverb prediction model, the dataset that contained definitions and synsets derived from a

list of adverbs and prepositions is carefully collected, playing a fundamental role in model training.

To enhance the precision of adverb prediction, the researchers incorporated the machine learning technique known as Latent Dirichlet Allocation (LDA), with specific application of the MLP (MultiLayer Perceptron) model. The researchers utilize the power of LDA to construct topic-based feature vectors for words, with a particular focus on adverbs. LDA is commonly used in NLP to discover hidden topics within a corpus of text. The process of generating these feature vectors comprised several key steps: first, LDA modeling was applied, wherein words were associated with specific topics to discover the underlying semantic patterns. Then, the *LDA_vector* method is introduced and designed to take a word as input and determine its LDA representation, representing the word as a vector of topic probabilities based on its contextual associations.

Additionally, the *Webster_LDA_vector* method is defined to extend this capability to adverbs not found in Wordnet but present in word embeddings, thereby broadening the scope of the LDA approach. Ultimately, the LDA-derived vectors obtained from these methods were integrated into the feature vectors for adverbs, providing a structured means to measure their similarity or categorization in the context of the discovered semantic topics. This feature-based analysis allowed for comprehensive comparisons with other word similarity measures, including spaCy and Wordnet-based metrics, enhancing our understanding of adverb similarities and categories.

In addition to the methodological approach, the researchers utilized the power of word embeddings. Word embeddings are a way to represent words as dense vectors in a continuous vector space, allowing us to capture relationships between words and how they fit into sentences. Within the scope of this study, the utilization of word embeddings offers several advantages. First, they help us measure how similar words are to each other, which is particularly useful for understanding adverbs in the context of other words. Second, when we encounter words that aren't in the dictionary (Wordnet) we're using in the code, word embeddings provide a smart solution by giving us vector representations for a wide range of words. Third, they enable us to understand the meaning of words within their context, making it easier to figure out what adverbs mean based on the words they're associated with. Fourth, when we're creating graphs that show how words relate to each other, word embeddings enhance these vectors with more information. This enrichment helps us better understand the roles of adverbs and other words in sentences. Lastly, the integration of word embeddings results in more accurate and detailed graphs, representing words and their connections in sentences, ultimately enhancing our overall understanding.

Now, with the understanding of how word embeddings enhance our analysis of word relationships, let's delve into the process of determining the functional type of a given sentence sequence. This process involves analyzing the structure and components of sentences to categorize them into different functional units. To do this, we consider a set of accepted functional unit types, which include Predicate, Subject, Direct Object, Indirect Object, Time, Place, Manner, Frequency, Degree, and Conjunction. This parsing process is the initial step in our study.

Having an input word sentence, $s = w_1, w_2, ..., w_l$. where symbol $w$ denotes a word inside the sentence. The set of accepted functional unit types is given by

$$T = \{\text{Predicate, Subject, Direct Object,Indirect Object,}$$
$$\text{Time,Place, Manner,Frequency, Degree, Conjunction}\} \tag{1}$$

To determine the functional type of a given sentence sequence, the following parsing processes are first:

1) The internal dictionary contains the list of frequent adverb words, like the phrase as soon as, in this case, the dictionary contains also the related functional type.

2) Label of the dependency parsing: $l_e$ This property is generated with the spacy parser as the label of the dependency edge from the generated dependency tree.

3) Wordnet-based Lin similarity ($l_l$): A score denoting how similar two word senses $(s_1, s_2)$ are, based on the Information Content (IC) of the Least Common Subsumer ($s_c$) most specific ancestor node) and that of the two input synsets:
$$l_l(s_1, s_2) = \frac{2 \cdot IC(s_c)}{IC(s_1) + IC(s_2)}$$

4) Wordnet-based path similarity ($l_p$): The path between the two synsets in the concept tree of the Wordnet.

5) Wordnet LDA similarity ($l_d$): We take the definition sections from the Wordnet database and calculate the topic similarities using the LDA method.

6) Webster LDA similarity ($l_w$): The definitions in Webster dictionary are used to calculate the topic similarities using the LDA method.

7) Spacy similarity ($l_s$): The similarity is based on the grammatical properties generated in the spacy NLP library.

The proposed framework also includes a dictionary which contains some selected words with the related unit types labels:

$$D = \{(w, T(w)\}$$

We divide this dictionary into two parts:

$$D = D_B \bigcup D_L$$

where $D_B$ is the set of baseline words, we use to determine the similarity positions of new query words. For a given query word $w_q$, the following local feature vectors are calculated:

$$\{l_e(w_q, w), l_l(w_q, w), l_p(w_q, w), l_d(w_q, w), l_w(w_q, w),$$
$$l_s(w_q, w) | w \epsilon D_B\}$$

Using these similarity measures, the generated similarity vectors are merged into a global feature vector

$$l(w_q)$$

These global feature vectors are used to predict the corresponding unit type label of $w_q$. For the prediction, an MLP neural network module (NN) is involved, where

$$NN(l(w_q))$$

outputs the predicted unit label.

```
Model: "sequential"
_____
 Layer (type)              Output Shape              Param #
=================================================================
 dense (Dense)             (None, 1100)              49500

 dropout (Dropout)         (None, 1100)              0

 dense_1 (Dense)           (None, 440)               484440

 dense_2 (Dense)           (None, 132)               58212

 dense_3 (Dense)           (None, 6)                 798

=================================================================
Total params: 592950 (2.26 MB)
Trainable params: 592950 (2.26 MB)
Non-trainable params: 0 (0.00 Byte)
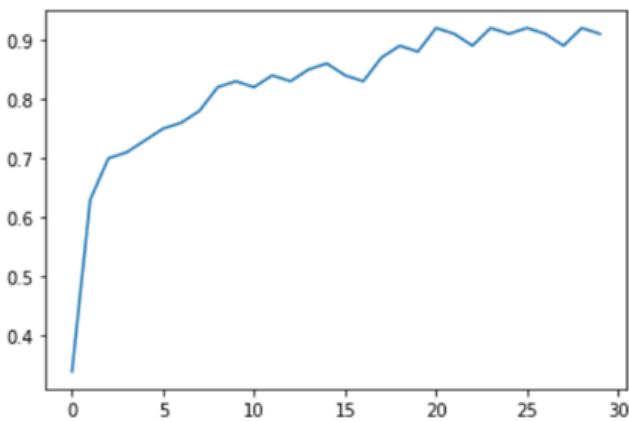```

Fig. 4. MLP architecture.



Fig. 5. Validation accuracy curve in the training process.

For the training of the MLP unit, the $D_L$ dataset is used as training and test dataset.

The MLP neural network unit under consideration comprises five layers, with one dedicated to model regularization (as depicted in Fig. 4). The trained MLP unit demonstrated a commendable average accuracy of 92% on the tested datasets.

Fig. 5 displays the validation accuracy curve during the training process of the proposed framework. The curve illustrates how the accuracy of the model evolves as it undergoes training iterations. It provides valuable insights into the model's performance and its ability to generalize to unseen data, showcasing the progress made during the training phase.

## V. SEMANTIC GRAPH INDUCTION

The process of automatically building a semantic graph from unstructured data, like textual documents or datasets, is known as semantic graph induction [26], [27]. It involves taking information from unstructured data, such as entities, concepts, and their relationships, and putting it in an organized manner. This procedure frequently depends on NLP and machine learning approaches to discover and link entities, infer relationships, and build the graph.

The term graphs refer to a common data format as well as a



Fig. 6. Application fields of knowledge graphs example. The famous zachary karate club network represents the friendship relationships between members of the karate club studied by Wayne W. Zachary from 1970 to 1972.

universal language for describing complicated systems. A common data structure and language for characterizing complex systems is called a graph. In its most basic form, a graph is just a set of objects or nodes, and the interactions (or edges) that exist between pairs of these nodes. For instance, we can utilize edges to signify the friendship between two individuals and utilize nodes to symbolize each person, effectively encoding a social network. This is illustrated in Fig. 6, featuring the renowned Zachary Karate Club Network.

An edge that connects two individuals if they socialize outside of the club. During Zachary's study, the club split into two factions centered around nodes 0 and 33 and Zachary was able to correctly predict which nodes would fall into each faction based on the graph structure [20]. Graphs do more than just provide an elegant theoretical framework, however. They offer a mathematical foundation that we can build upon to analyze, understand, and learn from real-world complex systems [20], [7].

Constructing large-scale semantic graphs from vast and diverse datasets is a significant challenge. Researchers are continually developing more efficient algorithms and technologies to handle big data [2]. Recently, word embeddings and entity embeddings have become effective in capturing semantic relationships, and the advancements in embedding techniques continue to improve graph construction [3]. Ensuring the completeness and accuracy of knowledge graphs is an ongoing challenge [4], [5], with methods for knowledge base completion and alignment being actively explored [6].

## VI. EXPERIMENTAL RESULTS

### A. Methodology

The dataset utilized for fine-tuning ChatGPT is meticulously curated from a wide array of linguistic sources, including academic texts in history and biology, and factual data about world events. This selection, aimed at capturing a rich variety of sentence structures, ensures exposure to complex and diverse linguistic patterns. Each sentence within this dataset is carefully labeled by linguistics experts to identify its functional components such as subjects, predicates, direct and indirect objects, as well as various types of adverbs like those indicating time, place, frequency, and manner. This detailed labeling is crucial for the accurate training of the model.

The size of the dataset was determined considering the resource-intensive nature of manual collection and analysis.

Fig. 7. The overall evaluation result of linguistic experts.

TABLE I. EFFICIENCY OF CHATGPT IN DEPENDENCY OF PROMPT SIZE

| Model | Prompt Size | Average Rating |
|---|---|---|
| OpenAI API | 5 | Poor |
| | 15 | Below Average |
| | 25 | Average |
| | 35 | Above Average |
| | 40 | Excellent |
| ChatGPT 3.5 Web Interface | – | Average |
| Hybrid Parser-based Sentence Parsing | – | Excellent |



Fig. 8. ChatGPT OpenAI and hybrid parser-based sentence parsing accuracy.

We assembled 160 sentences for the training dataset and 40 sentences for testing purposes. The semantic graph model employed in this study categorizes words and phrases from these sentences into their respective functional structures. This approach facilitates a comprehensive and nuanced understanding of sentence parsing, essential for the model's training and evaluation. The deliberate and diverse selection of sources ensures a well-rounded dataset, contributing to the effectiveness of the model in recognizing and interpreting a broad spectrum of linguistic elements.
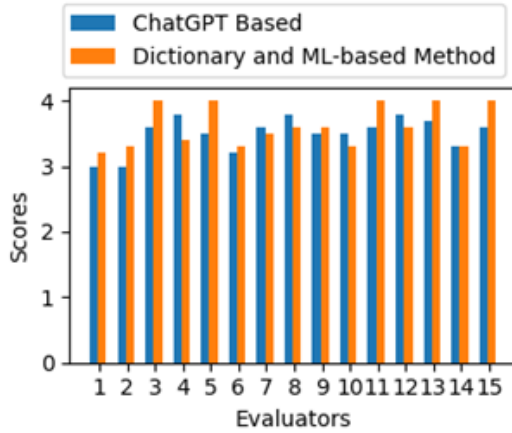
*B. Results and Analysis*

One significant limitation within this area of sentence parsing research is the absence of an automated performance evaluation system, which remains unimplemented. To assess the accuracy of the parsing, the researchers engaged the expertise of linguistic professionals, educators, and students. The survey encompassed five distinct rating categories: "Poor", "Below Average", "Average", "Above Average", and "Excellent". The researchers used similar test datasets for both approaches and make a comparative result analysis.

In the evaluation process, we used the ChatGPT efficiency for prompts of different lengths and complexity. The models evaluated in this study include the OpenAI API and ChatGPT 3.5 Web Interface, as well as a Hybrid Parser-based Method.

Fig. 7 provides an overview of the comprehensive evaluation results of 15 linguistic experts for both methods. The evaluation scores range from a minimum of 1.5 to a maximum of 4, showcasing the experts' assessments of the performance of these methods.

Efficiency, as reflected in the average quality rating of responses generated by these models, is a key measure. We explored prompt set sizes ranging from 5 to 40. Surprisingly, both the ChatGPT 3.5 Web Interface and the Hybrid Parser-based Sentence Parsing model consistently maintained an "excellent" response quality rating, irrespective of the prompt set size. This indicates their enduring efficiency across a spectrum of prompt set sizes. This table provides valuable insights into how different prompt set sizes impact ChatGPT model efficiency,

revealing noteworthy disparities in performance between the OpenAI API and other models.

Fig. 8 visually illustrates the influence of prompt set size on ChatGPT's sentence parsing performance, quantified by accuracy. Accuracy is determined by the ratio of correctly assigned sentences to the total assigned sentences. The OpenAI API employs five distinct prompts, each with varying numbers of sentence parsing templates: prompt one with 5 templates, prompt two with 15 templates, prompt three with 25 templates, prompt four with 35 templates, and prompt five with 40 templates. As seen in Table I, the accuracy of OpenAI models sees improvement as the number of templates within the prompts increases. In a separate experiment conducted with the ChatGPT 3.5 Web Interface, an accuracy score of 0.74 was achieved.

Table II presents accuracy values, indicating that the Hybrid Parser-based sentence parsing method exhibit a slight advantage over the ChatGPT-based model (acc_GPT = 0.85, acc_hybrid = 0.87). This evaluation scenario provides valuable insights into the performance and effectiveness of both approaches in sentence parsing.

This experiment underscores that while ChatGPT 3.5 is a recent and versatile language model capable of generating diverse and interesting results, it has limitations, particularly in domains like sentence parsing. The observed accuracy values

TABLE II. Efficiency of ChatGPT and Hybrid Parser-based Sentence Parsing Method

| Model | Prompt size | Accuracy |
|---|---|---|
| OpenAI API | 5 | 0.64 |
| | 15 | 0.67 |
| | 25 | 0.72 |
| | 35 | 0.77 |
| | 40 | **0.85** |
| ChatGPT 3.5 Web Interface | | 0.74 |
| Hybrid Parser-based Sentence Parsing | | **0.87** |



Fig. 9. Semantic graph generated by the proposed model for the sentence "Abebe reads a book deeply in the library each day after lunch".

strongly advocate for the effectiveness of the proposed Hybrid Parser based sentence parsing. This suggests that the proposed model may find broader applicability in sentence parsing tasks (see Fig. 9).

## VII. Discussion

This paper presents a novel approach to sentence parsing using ChatGPT, demonstrating significant potential in understanding and manipulating complex linguistic structures. We believe that the integration of LLMs like ChatGPT in sentence parsing tasks can revolutionize how we approach language understanding in AI. The model's ability to notice small details in language, from syntax to semantics, is particularly promising for applications in automated text summarization, sentiment analysis, and even in developing more advanced conversational AI.

However, we also recognize challenges, particularly in terms of computational demands and potential biases inherent in the training data. The scalability of such models in real-world applications remains a concern, especially considering the resource-intensive nature of their training and operation. It's crucial for future research to address these challenges, ensuring the responsible and efficient use of these powerful tools in various NLP applications.

## VIII. Conclusions and Future Work

In conclusion, the process of semantic graph construction stands as a cornerstone in the field of knowledge representation

and artificial intelligence, giving structured meaning upon the vast landscape of textual data. It draws its strength from an array of foundational technologies, encompassing NLP, dependency parsing, word embeddings, LDA, and the integration of ontologies and knowledge graphs. These technological underpinnings empower the creation of semantic graphs, spanning from entity recognition to intricate topic modeling. The infusion of ChatGPT's NLP capabilities further enriches this process, rendering it a dynamic and adaptable tool for semantic graph construction.
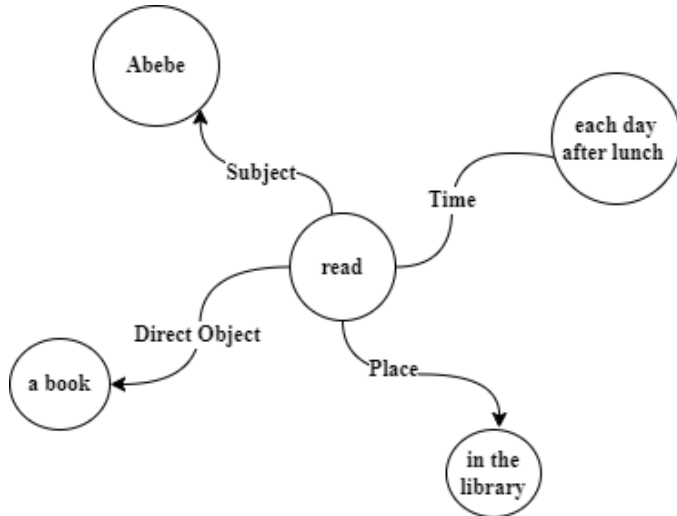
Our deliberate experimentation and meticulous evaluation have illuminated the comparative performance, applicability, and constraints of ChatGPT-based and Hybrid Parser based sentence parsing methods within the context of semantic graph construction. These findings not only contribute to the expanding reservoir of knowledge within the field of NLP but also offer invaluable insights to researchers, developers, and practitioners venturing into real-world applications. These applications include information retrieval, knowledge graph development, and automated question-answering systems, among others.

It's worth noting that the accuracy values indicate a slightly better performance of the hybrid parser-based sentence parsing method compared to the ChatGPT-based model $acc\_GPT = 0.85$, $acc\_hybrid = 0.87$. In this evaluation scenario, our test results provide comprehensive insights into the strengths and limitations of ChatGPT 3.5, particularly in the domain of English sentence parsing and language understanding tasks. This knowledge is instrumental in further enhancing the capabilities of ChatGPT for these specific tasks.

As we investigate into the future, ongoing efforts will focus on refining ChatGPT for improved performance in English sentence parsing, thus bridging the gap between language models and semantic graph construction. The integration of additional linguistic resources, enhanced fine-tuning techniques, and prompt engineering strategies will be explored to further empower ChatGPT in its role as a dynamic tool for language understanding and knowledge representation.

## References

[1] D. Rusu, B. Fortuna, M. Grobelnik, and D. Mladenić, "Semantic graphs derived from triplets application in document summarization," *Informatica*, vol. 33, no. 3, 2009.

[2] W. Hu, M. Fey, H. Ren, M. Nakata, Y. Dong, and J. Leskovec, "Ogb-lsc: A large-scale challenge for machine learning on graphs," *arXiv preprint arXiv:2103.09430*, 2021.

[3] P. Goyal and E. Ferrara, "Graph embedding techniques, applications, and performance: A survey," *Knowledge-Based Systems*, vol. 151, pp. 78–94, 2018. doi: 10.1016/j.knosys.2018.03.022.

[4] S. Ji, S. Pan, E. Cambria, P. Marttinen, and S. P. Yu, "A survey on knowledge graphs: Representation, acquisition, and applications," *IEEE transactions on neural networks and learning systems*, vol. 33, no. 2, pp. 494–514, 2021. doi: 10.1109/TNNLS.2021.3070843.

[5] B. Abu-Salih, "Domain-specific knowledge graphs: A survey," *Journal of Network and Computer Applications*, vol. 185, p. 103076, 2021. doi: 10.1016/j.jnca.2021.103076.

[6] H. Singh, P. Jain, S. Chakrabarti, et al., "Multilingual knowledge graph completion with joint relation and entity alignment," *arXiv preprint arXiv:2104.08804*, 2021.

[7] X. Zou, "A survey on application of knowledge graph," in *Journal of Physics: Conference Series*, vol. 1487, no. 1, 2020. doi: 10.1088/1742-6596/1487/1/012016.

[8] L. Shi, S. Li, X. Yang, J. Qi, G. Pan, B. Zhou, et al., "Semantic health knowledge graph: semantic integration of heterogeneous medical knowledge and services," *BioMed research international*, vol. 2017, 2017. doi: 10.1155/2017/2858423.

[9] H. Oh and R. Jain, "Detecting events of daily living using multimodal data," *arXiv preprint arXiv:1905.09402*, 2019.

[10] M. Rotmensch, Y. Halpern, A. Tlimat, S. Horng, and D. Sontag, "Learning a health knowledge graph from electronic medical records," *Scientific reports*, vol. 7, no. 1, p. 5994, 2017. doi: 10.1038/s41598-017-05778-z.

[11] J. Liu, Z. Lu, and W. Du, "Combining enterprise knowledge graph and news sentiment analysis for stock price prediction," 2019.

[12] S. A. Elnagdy, M. Qiu, and K. Gai, "Cyber incident classifications using ontology-based knowledge representation for cybersecurity insurance in financial industry," in *2016 IEEE 3rd International Conference on Cyber Security and Cloud Computing (CSCloud)*, 2016. doi: 10.1109/CSCloud.2016.45.

[13] F. Zablith, "Constructing social media links to formal learning: A knowledge Graph Approach," *Educational technology research and development*, vol. 70, no. 2, pp. 559–584, 2022. doi: 10.1007/s11423-022-10091-2.

[14] Y. Jia, Y. Qi, H. Shang, R. Jiang, and A. Li, "A practical approach to constructing a knowledge graph for cybersecurity," *Engineering*, vol. 4, no. 1, pp. 53–60, 2018. doi: 10.1016/j.eng.2018.01.004.

[15] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, et al., "Summary of ChatGPT-Related Research and Perspective Towards the Future of Large Language Models," *Meta-Radiology*, p. 100017, 2023. doi: 10.1016/j.metrad.2023.100017.

[16] M. Abdullah, A. Madain, and Y. Jararweh, "ChatGPT: Fundamentals, applications and social impacts," in *2022 Ninth International Conference on Social Networks Analysis, Management and Security (SNAMS)*, 2022. doi: 10.1109/SNAMS58071.2022.10062688.

[17] G. Vassiliou, N. Papadakis, and H. Kondylakis, "SummaryGPT: Leveraging ChatGPT for summarizing knowledge graphs" in *European Semantic Web Conference*, 2023. doi: 10.1007/978-3-031-43458-7_31.

[18] P. Pichappan, M. Krishnamurthy, and P. Vijayakumar, "Analysis of ChatGPT as a Question-Answering Tool,"

[19] H. B. de Barros Pereira, M. Grilo, I. de Sousa Fadigas, C. T. de Souza Junior, M. do Vale Cunha, R. S. F. Dantas Barreto, J. C. Andrade, and T. Henrique, "Systematic review of the 'semantic network' definitions," *Expert Systems with Applications*, pp. 118455, 2022. DOI: 10.1016/j.eswa.2022.118455.

[20] W. L. Hamilton, *Graph representation learning*. McGill University: Morgan & Claypool Publishers, 2020.

[21] Y. Chen, X. Ge, S. Yang, L. Hu, J. Li, and J. Zhang, "A Survey on Multimodal Knowledge Graphs: Construction, Completion and

[22] C. Peng, F. Xia, M. Naseriparsa, and F. Osborne, "Knowledge graphs: Opportunities and challenges," *Artificial Intelligence Review*, pp. 1–32, 2023. DOI: 10.1007/s10462-023-10465-9.

Applications," *Mathematics*, vol. 11, no. 8, pp. 1815, 2023. DOI: 10.3390/math11081815.

[23] A. Kamath and R. Das, "A survey on semantic parsing," *arXiv preprint arXiv:1812.00978*, 2018.

[24] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. N. Mendes, S. Hellmann, M. Morsey, P. Van Kleef, S. Auer, et al., "Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia," *Semantic web*, vol. 6, no. 2, pp. 167–195, 2015. DOI: 10.3233/SW-140134.

[25] Š. Čebirić, F. Goasdoué, H. Kondylakis, D. Kotzinos, I. Manolescu, G. Troullinou, and M. Zneika, "Summarizing semantic graphs: a survey," *The VLDB journal*, vol. 28, pp. 295–327, 2019. DOI: 10.1007/s00778-018-0528-3.

[26] P. Pham, L. T. T. Nguyen, W. Pedrycz, and B. Vo, "Deep learning, graph-based text representation and classification: a survey, perspectives and challenges," *Artificial Intelligence Review*, vol. 56, no. 6, pp. 4893–4927, 2023. DOI: 10.1007/s10462-022-10265-7.

[27] L. Kovacs, "Mapping Between Semantic Graphs and Sentences in Grammar Induction System," *BRAIN. Broad Research in Artificial Intelligence and Neuroscience*, vol. 1, pp. 100–112, 2010.

[28] H. Naveed, A. U. Khan, S. Qiu, M. Saqib, S. Anwar, M. Usman, N. Barnes, and A. Mian, "A comprehensive overview of large language models," *arXiv preprint arXiv:2307.06435*, 2023.

[29] E. Kotei and R. Thirunavukarasu, "A Systematic Review of Transformer-Based Pre-Trained Language Models through Self-Supervised Learning," *Information*, vol. 14, no. 3, pp. 187, 2023. DOI: 10.3390/info14030187.

[30] S. R. Bowman, "Eight things to know about large language models," *arXiv preprint arXiv:2304.00612*, 2023.

[31] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al., "Language models are few-shot learners," in *Advances in neural information processing systems*, vol. 33, 2020, pp. 1877–1901.

[32] C. D. Manning, "Human language understanding & reasoning," *Daedalus*, vol. 151, no. 2, pp. 127–138, 2022. DOI: 10.1162/daed_a_01905.

[33] B. Meskó, "Prompt Engineering as an Important Emerging Skill for Medical Professionals: Tutorial," *Journal of Medical Internet Research*, vol. 25, pp. e50638, 2023. DOI: 10.2196/50638.

[34] S. Vemprala, R. Bonatti, A. Bucker, and A. Kapoor, "ChatGPT for Robotics: Design Principles and Model Abilities," *Microsoft Technical Report*, February 2023. [Online]. Available: https://www.microsoft.com/en-us/research/ publication/

[35] P. P. Ray, "ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope," *Internet of Things and Cyber-Physical Systems*, vol. 3, pp. 121-154, 2023. ISSN: 2667-3452. DOI: 10.1016/j.iotcps.2023.04.003, [Online]. Available: https://www.sciencedirect.com/science/article/pii/S266734522300024X.

[36] G. Wohlgenannt, *Learning Ontology Relations by Combining Corpus-Based Techniques and Reasoning on Data from Semantic Web Sources*, Frankfurt am Main: Peter Lang International Academic Publishers, 2018.

[37] A. Harth, K. Hose, R. Schenkel, *Linked Data Management*, Boca Raton, FL, USA: CRC Press, 2014.

[38] F.-X. Desmarais, M. Gagnon, A. Zouaq, "Comparing a Rule-Based and a Machine Learning Approach for Semantic Analysis," in *Proceedings of the 6th International Conference on Advances in Semantic Processing*, Barcelona, Spain, 2012, pp. 103–108.

# Overview of Data Augmentation Techniques in Time Series Analysis

Ihababdelbasset ANNAKI, Mohammed RAHMOUNE, Mohammed BOURHALEB
Université Mohammed Premier, National School of Applied Sciences,
Laboratory of Research in Applied Sciences (LARSA),
Oujda, Morocco

*Abstract*—Time series data analysis is vital in numerous fields, driven by advancements in deep learning and machine learning. This paper presents a comprehensive overview of data augmentation techniques in time series analysis, with a specific focus on their applications within deep learning and machine learning. We commence with a systematic methodology for literature selection, curating 757 articles from prominent databases. Subsequent sections delve into various data augmentation techniques, encompassing traditional approaches like interpolation and advanced methods like Synthetic Data Generation, Generative Adversarial Networks (GANs), and Variational Autoencoders (VAEs). These techniques address complexities inherent in time series data. Moreover, we scrutinize limitations, including computational costs and overfitting risks. However, it's essential to note that our analysis does not end with limitations. We also comprehensively analyzed the advantages and applicability of the techniques under consideration. This holistic evaluation allows us to provide a balanced perspective. In summary, this overview illuminates data augmentation's role in time series analysis within deep and machine-learning contexts. It provides valuable insights for researchers and practitioners, advancing these fields and charting paths for future exploration.

*Keywords—Time series; data augmentation; machine learning; deep learning; synthetic data generation*

## I. INTRODUCTION

The concept of data augmentation has become indispensable in modern machine learning, serving as a key technique to enhance the diversity and volume of training data [1]. Its roots can be traced back to the early stages of machine learning, where the challenge of limited data first emerged. Augmentation techniques, through methods such as image rotation, flipping, or text paraphrasing, enable models to learn from a varied set of inputs, thereby increasing their generalization capabilities [2]. This is especially crucial in preventing overfitting, a common challenge in machine learning models trained on limited datasets [3].

Data augmentation transcends various learning paradigms, playing a significant role in both supervised and unsupervised learning contexts. In supervised learning, it addresses challenges like class imbalance and enriches small datasets, enhancing model accuracy and reliability [4]. In unsupervised learning, augmentation techniques help in extracting more robust features and patterns from unlabeled data, a vital aspect in domains such as natural language processing and computer vision [5]. The versatility of these techniques is also evident in their adaptability to different data types, including images, text, and audio [6], [7].

Time series data, with its sequential and often periodic nature, introduces unique augmentation challenges. Standard augmentation methods may not be directly applicable due to the temporal dependencies inherent in time series data. Techniques like time warping [8], window slicing, or injecting synthetic anomalies [9] are tailored to maintain these temporal relationships. Such methods have been shown to significantly improve the performance of models in various time series applications, from stock market predictions and weather forecasting to electrocardiogram analysis in healthcare [10].

Beyond improving model performance, data augmentation has broader impacts on the field of machine learning. It contributes to more efficient use of available data, reducing the need for extensive data collection, which can be costly and time-consuming. However, it also raises ethical considerations, particularly in ensuring that augmented data does not introduce or perpetuate biases. This is a critical aspect in applications involving human-centric data [11], [12], where fairness and representativeness are paramount.

This review provides a comprehensive analysis of data augmentation techniques with key contributions as follows:

- Holistic Overview: Showcases a wide array of data augmentation methods, presenting a broad perspective rather than focusing on a specific scope, thus providing a more inclusive understanding of the field.

- Comprehensive Analysis: Compared to earlier reviews, this approach stands out by offering a more thorough examination of data augmentation techniques across various machine learning and deep learning domains.

- Emphasis on Time Series Analysis: Particular attention is given to the applications and implications of these techniques in time series analysis, highlighting their relevance and utility in this specific area.

- Methodological Advancements: Covers the latest methodological advancements in data augmentation, providing insights into the evolving nature of these techniques.

- Real-World Applications and Cross-Domain Applicability: This review explores the practical applications and broad applicability of data augmentation techniques across various fields, highlighting their significant impact in real-world scenarios and their versatility in diverse contexts and domains.

Fig. 1. PRISMA (Preferred reporting items for systematic reviews and meta-analyses) [13].



Fig. 2. PRISMA (Preferred reporting items for systematic reviews and meta-analyses) for data augmentation in time series analysis.

- Pros and Effectiveness: Highlights the advantages and effectiveness of different data augmentation techniques, demonstrating their contribution to enhancing model performance and reliability.

- Limitations and Challenges: Addresses the limitations and challenges associated with data augmentation, offering a balanced view of their capabilities and constraints.

- Future Research Directions: Outlines potential future research directions, encouraging further exploration and development in the field of data augmentation.

The review is grounded in a systematic examination of a wide range of peer-reviewed literature, adhering to the PRISMA guidelines [13] (see Fig. 1).

The paper is structured to enhance comprehension, beginning with a methodology section that details the systematic approach to literature selection and analysis. Following that, subsequent sections delve into the specifics of data augmentation techniques, their applications in various real-world scenarios, their limitations and challenges, and conclude with a discussion on future research directions.

## II. RESEARCH METHODOLOGY FRAMEWORK

This overview was conducted adhering to the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) guidelines. While a formal pre-registered protocol was not established, the methodology was meticulously developed and documented prior to initiating the review, ensuring a structured and transparent approach.

The initial dataset for this review comprised a total of 757 peer-reviewed articles and preprints, identified using the specific research query "Data Augmentation" AND "Time Series" in major academic databases including preprints. This query was designed to capture studies published between 2019 and 2024 that specifically addressed the intersection of data augmentation techniques and time series analysis in the field of machine learning. To refine the dataset for relevance and accessibility, the articles were further screened based on language and access. The final selection criteria included articles published in English and available as open access. This filtering process narrowed the dataset down to 108 articles, ensuring a focused review of studies directly relevant to the core topic and broadly accessible to the research community. Articles that did not directly respond to the research query, and publications outside the specified time frame were excluded (see Fig. 2).

The selection process entailed a rigorous screening based on titles and abstracts to assess relevance, followed by a full-text review against the inclusion criteria. The study selection process was documented using a PRISMA flow diagram, which details the number of articles screened, assessed for eligibility, and included in the final review.

Data extraction was systematically conducted, focusing on extracting key information such as study objectives, methodologies, key findings, and specific techniques related to data augmentation. The extraction process was carried out by multiple reviewers to enhance accuracy, with any discrepancies resolved through consensus. A standardized data extraction template was employed to maintain consistency across all studies.

A bias assessment was performed using established criteria to evaluate the quality and reliability of each study. This assessment considered factors such as study design, methodology, data analysis, and reporting transparency.

Given the qualitative and diverse nature of the studies, a narrative synthesis approach was utilized. This involved identifying common themes, methodologies, and findings across the studies while considering the heterogeneity of the data and study designs.

The review was based on publicly available, published academic articles; therefore, it did not involve primary data collection or require ethical approval. The analysis was conducted with respect to the intellectual property of the original authors.

## III. STATISTICAL AND MACHINE LEARNING DATA AUGMENTATION TECHNIQUES

This section serves as an introduction to the diverse range of techniques encompassed by Statistical and Machine Learning Data Augmentation (see Fig. 3). It establishes the fundamental importance of data augmentation within the context of Time Series Analysis. By artificially expanding datasets and introducing variations, these techniques play a pivotal role in improving the robustness of models and the quality of insights drawn from time series data [14].

Within this subsection, we delve into the realm of statistical techniques used for data augmentation in time series analysis. Techniques such as Linear Interpolation enable the filling of gaps in data by estimating values between observed points, thus expanding datasets. Seasonal Decomposition separates time series into fundamental components, facilitating the generation of new samples by manipulating these constituent parts. Exponential Smoothing, on the other hand, focuses on forecasting future segments of time series data, effectively augmenting it with forward-looking information [15].

In this subsection, we shift our attention to Machine Learning-driven data augmentation approaches. Bootstrap Resampling enables the generation of multiple samples by randomly selecting data points with replacement, contributing to the diversification of datasets. K-Means Clustering partitions time series data into clusters based on similarity, allowing for the creation of new samples that exhibit different patterns [16]. Data Inpainting, a machine learning-based technique, aids in filling missing values by predicting them based on available data [17].

As we conclude this section, it's important to underscore the pivotal role that data augmentation plays in Time Series Analysis. By expanding datasets, improving data quality, and enabling the creation of synthetic samples, these techniques empower researchers and practitioners to extract more accurate insights from time series data [18]. The applicability of both statistical and machine learning methods underscores their relevance in a wide range of time series analysis tasks. Looking ahead, the continued development of data augmentation techniques promises to further advance the field, making it an area of ongoing interest and exploration (Table I).

## IV. DEEP LEARNING DATA AUGMENTATION TECHNIQUES

In this section, we explore advanced data augmentation techniques driven by Deep Learning models. These techniques are particularly effective in capturing complex patterns and dependencies within time series data, enabling the generation of high-quality synthetic samples.



Fig. 3. Statistical and machine learning data augmentation techniques.

TABLE I. SUMMARY OF DATA AUGMENTATION TECHNIQUES IN MACHINE LEARNING

| Technique | Description | Reference |
|---|---|---|
| Imputation Techniques | Explores the use of imputation methods for augmenting incomplete time series data, including techniques like Mean, Median, KNN-based imputation, Linear Regression, Miss Forest, and MICE to fill missing values. | [14], [15], [16], [17], [18] |
| Data Expansion Techniques | Discusses methods for augmenting datasets by expanding time series data, including techniques for urban expansion monitoring and forecasting using remote sensing data. | [19], [20], [21], [22], [23] |
| Time Series Transformation | Focuses on transforming time series data using machine learning techniques for augmentation, including methods for forecasting and analysis that enhance the richness of the dataset. | [24], [25], [26], [27], [28] |
| Statistical Models | Examines the use of statistical models for data augmentation in time series, comparing their performance with machine learning models in applications like heart failure event prediction. | [29], [30], [31], [32], [33] |
| Clustering and Similarity-Based Methods | Explores the application of clustering algorithms and similarity-based methods for augmenting datasets in machine learning, including use cases like customer segmentation and data analysis. | [34], [35], [36], [37], [38] |
| Data Sampling Techniques | Investigates various data sampling strategies for augmenting datasets in machine learning, especially for addressing imbalanced datasets in different domains. | [39], [40], [41], [42], [43] |

### A. Generative Models

*1) TimeGAN:* TimeGAN, a generative model designed for time series data, leverages a Generative Adversarial Network (GAN) framework to generate synthetic time series data that closely resembles the original data's statistical properties and dependencies [44], [45]. It comprises two main components: the generator and the discriminator. The generator aims to produce synthetic time series data, while the discriminator tries to distinguish between real and synthetic data [46], [47].

The loss function for TimeGAN is defined as:

$$\mathcal{L}_{\text{TimeGAN}} = \lambda \cdot \mathcal{L}_{\text{AdvD}} + (1 - \lambda) \cdot \mathcal{L}_{\text{AdvG}}$$

Here, $\mathcal{L}_{\text{AdvD}}$ represents the adversarial loss for the discriminator, $\mathcal{L}_{\text{AdvG}}$ is the adversarial loss for the generator, and $\lambda$ is a hyperparameter that balances the two losses [48].

*2) Variational Autoencoders (VAEs):* Variational Autoencoders (VAEs) are deep generative models that learn latent representations of time series data, used to generate new time series samples by sampling from the learned latent space [49], [50]. In a VAE, the encoder network maps the input time series data to a latent space where each point represents a potential data point, and the decoder network generates time series samples from points in the latent space [51], [52].

The loss function for VAEs consists of two terms: a reconstruction loss ($\mathcal{L}_{\text{rec}}$) that measures how well the generated data matches the original data and a regularization term ($\mathcal{L}_{\text{reg}}$) [53], [54]. This encourages the latent space to follow a predefined distribution, typically a Gaussian distribution. The loss is defined as:

$$\mathcal{L}_{\text{VAE}} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{reg}}$$

*3) Generative Adversarial Networks (GANs):* Generative Adversarial Networks (GANs) consist of a generator and a discriminator network that compete during training, and they are applied to generate synthetic time series data by training the generator to produce realistic samples. In a GAN, the generator aims to produce data that is indistinguishable from real data, while the discriminator tries to distinguish between real and generated data [55], [56].

The loss function for GANs is given by:

$$\mathcal{L}_{\text{GAN}} = E_{\text{real}}[\log(D(x))] + E_{\text{fake}}[\log(1 - D(G(z)))]$$

Here, $D(x)$ represents the discriminator's output for real data, $D(G(z))$ is the discriminator's output for generated data, and $z$ is a random noise vector [57], [58].

*4) LSTM Variational Autoencoders (LSTM-VAEs):* LSTM Variational Autoencoders (LSTM-VAEs) combine Long Short-Term Memory (LSTM) networks with VAEs for modeling and generating time series data, effectively capturing temporal dependencies [54], [49]. LSTM-VAEs consist of an encoder network that maps input time series data into a latent space and a decoder network that generates time series samples from points in the latent space [59], [60].

The loss function for LSTM-VAEs combines a reconstruction loss ($\mathcal{L}_{\text{rec}}$), similar to traditional VAEs, and a regularization term ($\mathcal{L}_{\text{reg}}$) that encourages the latent space to follow a predefined distribution [61]. The total loss is defined as:

$$\mathcal{L}_{\text{LSTM-VAE}} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{reg}}$$

*5) Temporal Generative Adversarial Networks (Temporal GANs):* Temporal Generative Adversarial Networks (Temporal GANs) specialize in generating time series data while considering the temporal nature of the data. Temporal GANs extend the traditional GAN framework to handle time series data. They use recurrent layers to capture temporal dependencies

and ensure that the generated data maintains the time sequence [55], [56].

The loss function for Temporal GANs is similar to the GAN loss but takes into account the sequential nature of the data, encouraging the generator to produce time-consistent samples.

*6) Wasserstein Generative Models:* Wasserstein Generative Models use the Wasserstein distance to measure data distribution similarity, aiming to create stable and high-quality synthetic time series data. The Wasserstein distance, also known as the Earth Mover's distance, quantifies the minimum amount of "work" required to transform one distribution into another. In the context of GANs, it provides a more stable and informative measure of the difference between real and generated data distributions [62], [63].

The loss function for Wasserstein GANs is defined as:

$$\mathcal{L}_{\text{WGAN}} = \sup_{\|D\|_L \leq 1} E_{\text{real}}[D(x)] - E_{\text{fake}}[D(G(z))]$$

Here, $D(x)$ represents the discriminator's output for real data, $D(G(z))$ is the discriminator's output for generated data, and $\|D\|_L \leq 1$ enforces a Lipschitz constraint on the discriminator.

*7) Recurrent Variational Autoencoders (RNN-VAE):* Recurrent Variational Autoencoders (RNN-VAE) employ recurrent neural networks (RNNs) and VAEs for modeling and generating sequential data, including time series.

RNN-VAEs incorporate RNN layers to handle sequential data and capture temporal dependencies. The encoder network maps input time series data to a latent space, and the decoder generates sequential data from points in the latent space.

The loss function for RNN-VAEs is similar to traditional VAEs, consisting of a reconstruction loss ($\mathcal{L}_{\text{rec}}$) and a regularization term ($\mathcal{L}_{\text{reg}}$) to encourage a predefined distribution in the latent space [64], [65], [66], [67], [68].

*8) Conditional Generative Models:* Conditional Generative Models allow for controlled generation based on specific conditions or input features.

In a conditional generative model, additional input information, known as conditions or context, is provided to the generator to influence the generation process. For example, conditions can include class labels or specific attributes that guide the generation of time series data.

The loss function for conditional generative models depends on the specific architecture and conditions used but typically involves both the reconstruction loss and a term related to the conditions used for generation [69], [70], [64], [71], [72], [73], [74], [75], [76], [77] (Table II).

*B. Sequence Modeling Techniques*

*1) Sequence-to-Sequence Models:* Sequence-to-sequence models are employed to generate new sequences based on the patterns learned from input sequences. They are widely used for time series data generation tasks [78].

TABLE II. GENERATIVE MODELS

| Technique | Description | References |
|---|---|---|
| TimeGAN | TimeGAN is a generative model designed for time series data. It leverages a Generative Adversarial Network (GAN) framework to generate synthetic time series data that closely resembles the original data's statistical properties and dependencies. | [44], [45], [46], [47], [48]. |
| Variational Autoencoders (VAEs) | Variational Autoencoders (VAEs) are deep generative models that can learn latent representations of time series data. They are used to generate new time series samples by sampling from the learned latent space. | [49], [50], [51], [52], [53], [54]. |
| Generative Adversarial Networks (GANs) | Generative Adversarial Networks (GANs) consist of a generator and a discriminator network that compete during training. They can be applied to generate synthetic time series data by training the generator to produce realistic samples. | [55], [56], [57], [58]. |
| LSTM Variational Autoencoders (LSTM-VAEs) | LSTM Variational Autoencoders (LSTM-VAEs) combine Long Short-Term Memory (LSTM) networks with VAEs for modeling and generating time series data. They are effective in capturing temporal dependencies. | [54], [49], [59], [60], [61] |
| Temporal Generative Adversarial Networks (Temporal GANs) | Temporal GANs are specialized GANs for generating time series data. They consider the temporal nature of the data during the generation process. | [55], [56] |
| Wasserstein Generative Models | Wasserstein Generative Models use the Wasserstein distance to measure the similarity between real and generated data distributions. They aim to create more stable and high-quality synthetic time series data. | [62], [63] |
| Recurrent Variational Autoencoders (RNN-VAE) | Recurrent Variational Autoencoders (RNN-VAE) employ recurrent neural networks (RNNs) and VAEs to model and generate sequential data, including time series. | [64], [65], [66], [67], [68] |
| Conditional Generative Models | Conditional generative models generate data samples based on specific conditions or input features, allowing for the controlled generation of time series data. | [69], [70], [64], [71], [72], [73], [74], [75], [76], [77] |

*2) Data Augmentation through Noise Addition:* Data Augmentation through Noise Addition involves injecting controlled noise into the time series data to generate variations and enhance the training dataset. This approach can be represented as follows: Given an original time series $\mathbf{X} = [x_1, x_2, \ldots, x_T]$, where $x_t$ represents the value at time $t$, and a noise signal $\mathbf{N} = [n_1, n_2, \ldots, n_T]$, where $n_t$ is sampled from a predefined noise distribution, the augmented time series is obtained as $\mathbf{X}_{\text{aug}} = \mathbf{X} + \mathbf{N}$ [79].

*3) Transformer Models:* Transformer Models, known for their effectiveness in sequence modeling tasks, can be used to generate time series data by modeling long-range dependencies. The Transformer architecture includes self-attention mechanisms, which can capture relationships between distant time steps [80].

*4) Temporal Convolutional Networks:* Temporal Convolutional Networks (TCNs) utilize convolutional layers to capture temporal patterns in time series data and generate new sequences. A 1D convolutional layer with kernel size $K$ is used to capture local patterns in TCNs [82].

## V. REAL-WORLD APPLICATIONS AND USE CASES OF DATA AUGMENTATION IN TIME SERIES ANALYSIS

Data augmentation techniques have found invaluable applications in various real-world scenarios within the field of time series analysis. These methods are employed to tackle specific challenges, enhance predictive models, and enable more accurate forecasts across diverse domains.

In the realm of finance, data augmentation plays a pivotal role in generating synthetic financial time series data. This synthetic data supplements genuine financial records and is particularly useful in training predictive models for stock market analysis and portfolio management. For instance, the effectiveness of LSTM-GAN in generating synthetic time series data, achieving a close resemblance to real data with similar silhouette scores and low Mean Squared Error (MSE) and Root Mean Squared Error (RMSE) values, was demonstrated by Chen et al. [81]. Furthermore, S. Crepey et al. [82] proposed an approach to improve anomaly detection in financial time series, showing that value-at-risk estimation errors are reduced when using the proposed model. By introducing simulated market conditions and variations, data augmentation contributes to the development of robust financial models."

In the healthcare and medical research sectors, privacy regulations and limited access to patient data can pose significant hurdles. Data augmentation techniques come to the rescue by creating synthetic patient time series data. Yang et al. developed TS-GAN, a Time-series GAN based on LSTM networks, to augment sensor-based health data in healthcare. This approach significantly enhances the performance of classification models, achieving classification accuracies of 97.50% on ECG_200, 94.12% on NonInvasiveFatalECG_Thorax1, and 98.12% on mHealth datasets [83]. Furthermore, the improvement of SAX representation for time series using wavelet packet decomposition and FastDTW by Guo et al. [84] has the highest classification accuracy in 11 of 20 datasets. These artificial datasets empower the development of predictive models for disease diagnosis, patient monitoring, and drug discovery, all while safeguarding patient privacy and complying with data regulations.

Within the manufacturing and industrial domains, data augmentation strategies involve generating synthetic sensor data and introducing anomalies into existing datasets. This augmented data enhances the resilience of predictive maintenance

models, resulting in improved equipment uptime and operational efficiency. For instance, the application of simulation-based data augmentation for the quality inspection of structural adhesive with deep learning improved the performance of models in a scarce manufacturing data context with imbalanced training sets by 3.1% (mAP@0.50) [85]. Additionally, strategic data augmentation with CTGAN for smart manufacturing significantly enhanced machine learning predictions of paper breaks in pulp-and-paper production. The models' detection of machine breaks improved by over 30% for Decision Trees, 20% for Random Forest, and nearly 90% for Logistic Regression [86]. These advancements underscore data augmentation as a critical component of predictive maintenance and process optimization in industrial settings.

The energy and utilities industry leverages data augmentation to simulate energy consumption and production variations. This synthetic data aids in forecasting energy demand, optimizing grid operations, and ensuring a stable energy supply [87]Data augmentation appears to have significantly improved the forecasting accuracy in both the univariable and multivariable models. This is evident from the lower RMSE and MAPE values across all regions when comparing the augmented columns to their non-augmented counterparts. For instance, looking at the Busan region: The RMSE for the univariable model without augmentation is 0.2345, and with augmentation is 0.0853, showing a marked improvement. The RMSE for the multivariable model without augmentation is 0.1722 and with augmentation is 0.0132, which is a significant decrease. Augmented time series data contributes to effective resource management and reduced disruptions in the energy sector.

Environmental monitoring relies on data augmentation to replicate variations in environmental factors and weather conditions. Specifically, in the case of crack detection in AGR and CFD data as discussed by Branikas et al. in 2023 [88], the augmentation demonstrates a noticeable enhancement in recall and F1 score when applying a small pixel relaxation radius. Importantly, this dataset was not annotated using specialized tools or assessed by human experts. These synthetic time series datasets complement real-world observations, thereby contributing to more precise weather predictions, air quality assessments, and early detection of natural disasters. Augmentation remains a vital component in proactive environmental management and disaster preparedness.

In summary, data augmentation techniques are indispensable in time series analysis across a wide array of real-world applications and use cases. Whether in finance, healthcare, manufacturing, energy, environmental monitoring, or IoT, these methods empower the development of predictive models, improve operational efficiency, and support critical decision-making processes.

## VI. Challenges and Limitations of Time Series Data Augmentation Techniques

While time series data augmentation techniques offer significant advantages in various applications, they are not without their challenges and limitations. Understanding these constraints is essential for making informed decisions when employing these methods.

### A. Preservation of Temporal Dependencies

One of the primary challenges in time series data augmentation is the preservation of temporal dependencies. Many real-world time series exhibit complex dependencies and patterns over time. Data augmentation techniques must ensure that synthetic data maintains these dependencies accurately [89]. In cases where temporal structures are not adequately preserved, the performance of predictive models may degrade [90].

### B. Quality of Synthetic Data

The quality of synthetic data generated through augmentation techniques is a critical concern [91]. The synthetic data should closely resemble real-world observations to ensure that predictive models trained on augmented data generalize effectively. Poorly generated synthetic data can introduce biases and inaccuracies, leading to unreliable model outcomes [92].

### C. Generalization to Unseen Scenarios

Data augmentation should enable predictive models to generalize well to unseen scenarios [93]. However, there is a risk that the augmented data may be too tailored to specific training conditions, limiting the model's ability to handle novel situations [94]. Striking a balance between augmentation and maintaining generalization capabilities is a challenging task.

### D. Data Privacy and Ethical Considerations

In certain domains, such as healthcare and finance, data privacy and ethical concerns pose limitations on the use of data augmentation techniques [95]. Creating synthetic patient or financial data must adhere to strict privacy regulations and ethical guidelines, which can be a complex and resource-intensive process.

### E. Computational Complexity

Some advanced data augmentation techniques, particularly those involving generative models can be computationally intensive and time-consuming [96]. The computational complexity of generating large volumes of synthetic data may limit the scalability of augmentation methods.

### F. Availability of Domain-Specific Augmentation Tools

The availability of domain-specific data augmentation tools and expertise can be limited [89]. Applying augmentation techniques effectively often requires domain knowledge and specialized software, which may not be readily accessible in all applications.

### G. Evaluation and Validation

Evaluating the effectiveness of data augmentation methods and validating the performance of predictive models trained on augmented data can be challenging [90]. Developing appropriate evaluation metrics and conducting rigorous testing are essential but can be time and resource-intensive.

In conclusion, while time series data augmentation techniques offer numerous advantages, they also come with challenges and limitations that must be carefully considered. Addressing these limitations and understanding the constraints of each technique is crucial to ensure the successful application of data augmentation in time series analysis.

## VII. Comprehensive Analysis of Data Augmentation Techniques: Advantages, Limitations, and Applicability

In the evolving landscape of machine learning and data science, data augmentation techniques play a pivotal role in enhancing model performance and reliability. These techniques are instrumental in addressing challenges such as data scarcity, imbalanced datasets, and overfitting. This section provides a thorough analysis of various data augmentation techniques, exploring their advantages, limitations, and ideal use cases.

Table III presents a comprehensive examination of both traditional and advanced data augmentation techniques, encompassing methods ranging from Imputation Techniques to cutting-edge approaches like TimeGAN, Variational Autoencoders (VAEs), and Transformer Models. The table assesses each technique's effectiveness, potential drawbacks, and the scenarios where they are most beneficial. This includes an exploration of traditional data augmentation methods as well as advanced generative models and sequence modeling techniques.

These comprehensive tables serve as a guide for researchers and practitioners to select the most appropriate data augmentation strategies, tailored to the specific needs and constraints of their machine-learning projects.

## VIII. Conclusion

Time series analysis is a fundamental component of various domains, including finance, healthcare, environmental science, and more. The success of predictive models in these fields often hinges on the availability of diverse and high-quality time series data. However, obtaining such data can be challenging due to limited samples, data privacy concerns, or resource constraints. To address these challenges, data augmentation techniques have emerged as valuable tools in the time series analyst's toolkit.

In this paper, we provided an in-depth overview of data augmentation techniques in time series analysis. We explored various categories of augmentation methods, from statistical techniques to machine learning and deep learning approaches. Each category offers unique advantages and is applicable to different use cases.

Statistical techniques, such as linear interpolation, seasonal decomposition, and rolling window aggregation, provide simple and interpretable ways to augment time series data. Machine learning methods, like bootstrapping, semi-supervised learning, and time series embeddings, offer more sophisticated approaches for generating synthetic data. Deep learning techniques, including GANs, VAEs, and sequence-to-sequence models, push the boundaries of data augmentation by creating highly realistic and complex synthetic time series.

We delved into the mathematical foundations and practical applications of these techniques, showcasing their utility in tasks such as forecasting, anomaly detection, and trend analysis. Moreover, we discussed real-world use cases in finance, healthcare, and environmental monitoring, highlighting the impact of data augmentation on improving model performance and decision-making.

However, it is crucial to acknowledge that data augmentation in time series analysis is not without its challenges and limitations. Preserving temporal dependencies, ensuring data quality, and addressing computational complexity are ongoing concerns. Ethical considerations and domain-specific requirements further complicate the adoption of these techniques.

In conclusion, data augmentation techniques in time series analysis offer a promising avenue to tackle data scarcity and enhance the capabilities of predictive models. Researchers and practitioners should carefully assess the suitability of these techniques for their specific applications while being mindful of their limitations. The ever-evolving landscape of data augmentation continues to expand, opening doors to new possibilities in time series analysis and beyond.

## IX. Future Research Directions

As data augmentation techniques in time series analysis continue to evolve and gain prominence, several promising avenues for future research emerge. These directions are expected to shape the field and address existing challenges while opening up new possibilities for innovation. In this section, we outline some key areas for future exploration:

- One critical area of research is the development of data augmentation methods that better preserve temporal dependencies within time series data [97].

- As data augmentation becomes more prevalent, ethical considerations surrounding the generation and use of synthetic data warrant careful examination [98].

- Expanding the applicability of data augmentation techniques to cross-domain scenarios is an exciting direction for research [99].

- Hybrid data augmentation approaches that combine statistical, machine learning, and deep learning methods offer a promising avenue for exploration [100].

- Integrating data augmentation into automated machine learning (AutoML) pipelines can streamline the model development process [101].

- Interpretable and explainable data augmentation methods are essential for building trust in augmented data and the models trained on them [102].

- Establishing standardized benchmark datasets and evaluation metrics for assessing the quality and performance of data augmentation techniques is crucial [103].

- Efforts to design resource-efficient data augmentation techniques, especially for scenarios with limited computational resources, are essential [104].

In summary, the field of data augmentation in time series analysis offers abundant opportunities for future research and innovation. Researchers and practitioners can delve into areas such as preserving temporal dependencies, addressing ethical concerns, exploring cross-domain applications, and seamlessly integrating data augmentation into AutoML processes. As data augmentation remains pivotal in enhancing time series analysis, staying at the forefront of these research directions becomes imperative to unleash its full potential.

TABLE III. ADVANTAGES, LIMITATIONS, AND APPLICABILITY OF DATA AUGMENTATION TECHNIQUES IN MACHINE AND DEEP LEARNING

| Technique | Advantages | Limitations | Applicability |
|---|---|---|---|
| Imputation Techniques | - Can effectively handle missing data, improving dataset completeness.<br>- Offers a variety of methods suitable for different data types and patterns. | - Risk of introducing bias or inaccuracies, especially if the imputation model doesn't align well with the data's nature.<br>- Might oversimplify complex data relationships. | - Best used when dealing with datasets having missing values, especially in cases where the data is crucial and cannot be discarded. |
| Data Expansion Techniques | - Allows for the creation of larger and more diverse datasets.<br>- Particularly useful in fields like remote sensing where data can be scarce. | - Expanded data might not always represent real-world scenarios accurately.<br>- Risk of introducing artificial patterns not present in the original dataset. | - Ideal for situations where the available dataset is too small or lacks diversity, such as in certain types of research or specialized applications. |
| Time Series Transformation | - Enhances the diversity and richness of data, leading to potentially better model performance.<br>- Useful for both forecasting and deeper data analysis. | - Transformation techniques can distort the original time series properties.<br>- Requires careful selection to ensure relevance and accuracy. | - Suitable for time series forecasting, especially when the goal is to reveal hidden patterns or to adapt data to specific analytical needs. |
| Statistical Models | - Provides a more traditional and often simpler approach to data augmentation.<br>- Good for understanding underlying data distributions. | - May not capture complex nonlinear relationships as effectively as more advanced machine learning models.<br>- Limited flexibility in handling diverse data types. | - Recommended for scenarios where a straightforward, interpretable approach is needed, particularly in fields with well-understood data distributions. |
| Clustering and Similarity-Based Methods | - Useful for discovering natural groupings and patterns in data.<br>- Can improve data organization and segmentation. | - Performance is heavily dependent on the choice of similarity measures.<br>- Can be sensitive to outliers and noise in the data. | - Best applied in data segmentation, customer profiling, or any scenario requiring the identification of inherent groupings in the data. |
| Data Sampling Techniques | - Effective in addressing imbalanced datasets, and enhancing model training.<br>- Various strategies available to suit different data scenarios. | - Risks include overfitting, underfitting, or introducing sampling bias.<br>- This may lead to loss of important information if not carefully implemented. | - Particularly useful in cases of imbalanced datasets, such as in fraud detection or rare event prediction, where certain classes are underrepresented. |
| TimeGAN | - Excellent for capturing temporal dynamics in time series.<br>- Generates data that closely resembles real statistical properties. | - Computationally intensive.<br>- Requires large amounts of training data for accuracy. | - Ideal for scenarios where authentic-like time series data generation is needed, such as financial market analysis. |
| Variational Autoencoders (VAEs) | - Good at learning complex distributions.<br>- Capable of generating diverse data samples. | - Can struggle with generating high-quality reconstructions.<br>- Somewhat complex to train and tune. | - Suitable for tasks requiring the generation of new samples from complex data distributions, like image or speech synthesis. |
| Generative Adversarial Networks (GANs) | - Can produce highly realistic synthetic data.<br>- Versatile for various data types. | - Training can be unstable.<br>- Prone to mode collapse. | - Best for applications where realistic data generation is crucial, such as art creation or data augmentation. |
| LSTM Variational Autoencoders (LSTM-VAEs) | - Effective in modeling time dependencies.<br>- Combines LSTM's sequence handling with VAE's generative capabilities. | - Risk of overfitting on smaller datasets.<br>- Complex model architecture. | - Useful in sequential data applications like anomaly detection in time series. |
| Temporal Generative Adversarial Networks (Temporal GANs) | - Specifically designed for time series data.<br>- Addresses temporal aspects effectively. | - Can be computationally demanding.<br>- Requires careful tuning and training. | - Ideal for generating time-dependent synthetic data, such as in healthcare or stock market prediction. |
| Wasserstein Generative Models | - Offers more stable training than traditional GANs.<br>- Better at handling data distribution. | - More challenging to implement.<br>- Can be computationally more intensive. | - Recommended for scenarios where stable training of generative models is a priority, like in large-scale data generation. |
| Recurrent Variational Autoencoders (RNN-VAE) | - Good for sequential data representation.<br>- Combines RNN's temporal modeling with VAE's generative properties. | - Training can be time-consuming.<br>- Susceptible to vanishing gradients problem. | - Suitable for generating complex time series or sequential data, such as in natural language processing. |
| Conditional Generative Models | - Allows control over generated data features.<br>- Highly versatile in data generation. | - Requires additional conditioning data.<br>- Increased model complexity. | - Best used when specific conditions or features need to be included in the generated data, like in targeted marketing campaigns. |
| Sequence-to-Sequence Models | - Effective for generating sequences based on learned patterns.<br>- Widely applicable in time series generation. | - Requires large amounts of data for accuracy.<br>- Can be complex to tune and optimize. | - Ideal for applications like machine translation, speech recognition, and time series forecasting. |
| Data Augmentation through Noise Addition | - a simple and effective way to create data variations.<br>- Enhances the robustness of models. | - Risk of distorting the original data too much.<br>- Noise parameters need to be carefully chosen. | - Useful in scenarios where minor variations in the dataset can lead to significant improvements, such as in image or signal processing. |
| Transformer Models | - Excellent at capturing long-range dependencies.<br>- Self-attention mechanism provides dynamic focus. | - Can be resource-intensive.<br>- Requires significant amounts of training data. | - Suitable for complex sequence modeling tasks like natural language understanding and time series analysis. |
| Temporal Convolutional Networks (TCNs) | - Effective in capturing local and global temporal patterns.<br>- Efficient in terms of computational resources. | - May miss intricate long-term dependencies.<br>- Architecture needs a careful design for specific tasks. | - Recommended for tasks like audio synthesis and real-time anomaly detection in time series data. |

However, it's crucial to acknowledge certain limitations in our comprehensive overview. Our scope may not cover all existing techniques, and the diverse nature of time series data, along with the choice of evaluation metrics, may limit generalizability. Overfitting risks, the ever-evolving research landscape, interdisciplinary variations, and data accessibility issues are additional factors that deserve attention. Despite these challenges, our goal was to furnish a balanced and informative overview, serving as a valuable guide for both researchers and practitioners in the field.

REFERENCES

[1] M. Rahman, M. Rivolta, F. Badilini, R. Sassi, "A Systematic Survey of Data Augmentation of ECG Signals for AI Applications," *Sensors*, vol. 23, no. 11, 2023. doi:10.3390/s23115237

[2] N. A. Andriyanov, D. Andriyanov, "The using of data augmentation in machine learning in image processing tasks in the face of data scarcity," *Journal of Physics: Conference Series*, vol. 1661, no. 1, 2020. doi:10.1088/1742-6596/1661/1/012018

[3] S. Aleem, T. Kumar, S. Little, M. Bendechache, R. Brennan, K. McGuinness, "Random Data Augmentation based Enhancement: A Generalized Enhancement Approach for Medical Datasets," *Frontiers in Medicine*, 2022. doi:10.56541/fumf3414

[4] C. M. Burlacu, A. Burlacu, M. Praisler, C. Paraschiv, "Harnessing Deep Convolutional Neural Networks Detecting Synthetic Cannabinoids: A Hybrid Learning Strategy for Handling Class Imbalances in Limited Datasets," *Inventions*, vol. 8, no. 5, 2023. doi:10.3390/inventions8050129

[5] C.-Y. Hsu, P.-Y. Chen, S. Lu, S. Liu, C.-M. Yu, "Adversarial Examples Can Be Effective Data Augmentation for Unsupervised Machine Learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 6, 2021. doi:10.1609/aaai.v36i6.20650

[6] Q. Xie, Z. Dai, E. Hovy, M.-T. Luong, Q. V. Le, "Unsupervised Data Augmentation for Consistency Training," *arXiv*, 2019. [Online]. Available: https://arxiv.org/abs/1904.12848

[7] J. Yoo, T. Zhao, L. Akoglu, "Understanding the Effect of Data Augmentation in Self-supervised Anomaly Detection," *arXiv*, 2022. doi:10.48550/arXiv.2208.07734

[8] B. K. Iwana, S. Uchida, "Time Series Data Augmentation for Neural Networks by Time Warping with a Discriminative Teacher," *IEEE International Conference on Pattern Recognition*, 2020. doi:10.1109/ICPR48806.2021.9412812

[9] A. Aboussalah, M. Kwon, R. G. Patel, C. Chi, C.-G. Lee, "Don't overfit the history - Recursive time series data augmentation," *arXiv*, 2022. doi:10.48550/arXiv.2207.02891

[10] X. Yang, Z. Zhang, X. Cui, R.-y. Cui, "A Time Series Data Augmentation Method Based on Dynamic Time Warping," *IEEE Conference*, 2021. doi:10.1109/CCAI50917.2021.9447507

[11] I. Pastaltzidis, N. Dimitriou, K. Quezada-Tavárez, S. Aidinlis, T. Marquenie, A. Gurzawska, D. Tzovaras, "Data augmentation for fairness-aware machine learning: Preventing algorithmic bias in law enforcement systems," *ACM Conference*, 2022. doi:10.1145/3531146.3534644

[12] J. Yuan, R. Tang, X. Jiang, X. Hu, "LLM for Patient-Trial Matching: Privacy-Aware Data Augmentation Towards Better Performance and Generalizability," *arXiv*, 2023. doi:10.48550/arXiv.2303.16756

[13] M. Zuccon, E. Topino, A. Musetti, A. Gori, "Psychodynamic Therapies for the Treatment of Substance Addictions: A PRISMA Meta-Analysis," *Journal of Personalized Medicine*, vol. 13, no. 10, 2023. doi:10.3390/jpm13101469

[14] E. A. Christobel, "Imputation Techniques in Machine Learning – A Survey," *International Journal of Recent Technology and Engineering*, vol. 11, no. 10, 2023. doi:10.17762/ijritcc.v11i10.8662

[15] H. Wang et al., "Application of machine learning missing data imputation techniques in clinical decision making," *BMC Medical Informatics and Decision Making*, vol. 22, no. 1, 2022. doi:10.1186/s12911-022-01752-6

[16] A. R. Ismail, N. Z. Abidin, and M. Maen, "Systematic Review on Missing Data Imputation Techniques with Machine Learning Algorithms for Healthcare," *Journal of Robotics and Control (JRC)*, vol. 3, no. 2, 2022. doi:10.18196/jrc.v3i2.13133

[17] V. P. C. Magboo et al., "Imputation Techniques and Recursive Feature Elimination in Machine Learning Applied to Type II Diabetes Classification," *ACM International Conference Proceeding Series*, 2021. doi:10.1145/3508259.3508288

[18] T. Thomas and E. Rajabi, "A systematic review of machine learning-based missing value imputation techniques," *Data Technologies and Applications*, vol. 55, no. 3, 2021. doi:10.1108/DTA-12-2020-0298

[19] E. Mostafa et al., "Monitoring and Forecasting of Urban Expansion Using Machine Learning-Based Techniques and Remotely Sensed Data: A Case Study of Gharbia Governorate, Egypt," *Remote Sensing*, vol. 13, no. 22, 2021. doi:10.3390/rs13224498

[20] A. Agnihotri et al., "Role of data mining and machine learning techniques in medical imaging," *International Journal of Advanced Intelligence Paradigms*, vol. 16, no. 1/2, 2020. doi:10.1504/IJAIP.2018.10017086

[21] S. Fha et al., "Development of an Efficient Method to Detect Mixed Social Media Data with Tamil-English Code Using Machine Learning Techniques," *ACM International Conference Proceeding Series*, 2022. doi:10.1145/3563775

[22] A. Nafees et al., "Forecasting the Mechanical Properties of Plastic Concrete Employing Experimental Data Using Machine Learning Algorithms: DT, MLPNN, SVM, and RF," *Polymers*, vol. 14, no. 8, 2022. doi:10.3390/polym14081583

[23] E. Aharoni et al., "HE-PEx: Efficient Machine Learning under Homomorphic Encryption using Pruning, Permutation and Expansion," *arXiv preprint arXiv:2207.03384*, 2022. doi:10.48550/arXiv.2207.03384

[24] C. Kubik et al., "Knowledge discovery from time series in engineering applications using machine learning techniques," *Journal of Manufacturing Science and Engineering, Transactions of the ASME*, vol. 144, no. 3, 2022. doi:10.1115/1.4054158

[25] D. Salwala et al., "Distributed Incremental Machine Learning for Big Time Series Data," *IEEE International Conference on Big Data (Big Data)*, 2022. doi:10.1109/BigData55660.2022.10020361

[26] F. Ott et al., "Domain Adaptation for Time-Series Classification to Mitigate Covariate Shift," *ACM International Conference Proceeding Series*, 2022. doi:10.1145/3503161.3548167

[27] "The composition of time-series images and using the technique SMOTE ENN for balancing datasets in land use/cover mapping," *Applied Mathematics and Sciences: An International Journal*, vol. 27, no. 2, 2022. doi:10.46544/ams.v27i2.05

[28] S. Kuili et al., "A holistic machine learning approach to identify performance anomalies in enterprise WiFi deployments," *Proceedings of SPIE - The International Society for Optical Engineering*, vol. 12118, 2022. doi:10.1117/12.2621087

[29] Z. Sun et al., "Comparing Machine Learning Models and Statistical Models for Predicting Heart Failure Events: A Systematic Review and Meta-Analysis," *Frontiers in Cardiovascular Medicine*, vol. 9, 2022. doi:10.3389/fcvm.2022.812276

[30] Y. Li et al., "Consistency of variety of machine learning and statistical models in predicting clinical risks of individual patients: longitudinal cohort study using cardiovascular disease as exemplar," *BMJ (Clinical research ed.)*, vol. 371, 2020. doi:10.1136/bmj.m3919

[31] D. Zhou et al., "Integration of machine learning and statistical models for crash frequency modeling," *Journal of Intelligent Transportation Systems: Technology, Planning, and Operations*, 2022. doi:10.1080/19427867.2022.2158257

[32] M. S. Jaafarzadeh et al., "Groundwater recharge potential zonation using an ensemble of machine learning and bivariate statistical models," *Scientific Reports*, vol. 11, no. 1, 2021. doi:10.1038/s41598-021-85205-6

[33] X. Dastile et al., "Statistical and machine learning models in credit scoring: A systematic literature survey," *Applied Soft Computing*, vol. 96, 2020. doi:10.1016/j.asoc.2020.106263

[34] L. Kwuida and D. Ignatov, "On Interpretability and Similarity in Concept-Based Machine Learning," *Advances in Intelligent Systems and Computing*, vol. 1260, 2021. doi:10.1007/978-3-030-72610-2_3

[35] K. Shahina and T. P. Kumar, "Similarity-based clustering and data aggregation with independent component analysis in wireless sensor networks," *Transactions on Emerging Telecommunications Technologies*, 2022. doi:10.1002/ett.4462

[36] H. Zhu et al., "Assessment of the Generalization Abilities of Machine-Learning Scoring Functions for Structure-Based Virtual Screening," *Journal of Chemical Information and Modeling*, 2022. doi:10.1021/acs.jcim.2c01149

[37] D. M. S. Shakoor et al., "A Machine Learning Recommender System Based on Collaborative Filtering Using Gaussian Mixture Model Clustering," *Authorea Preprints*, 2020. doi:10.22541/au.160897179.93005705/v1

[38] K. Polat, "Similarity-based attribute weighting methods via clustering algorithms in the classification of imbalanced medical datasets," *Neural Computing and Applications*, 2018. doi:10.1007/s00521-018-3471-8

[39] M. Imani and H. Arabnia, "Hyperparameter Optimization and Combined Data Sampling Techniques in Machine Learning for Customer Churn Prediction: A Comparative Analysis," *Preprints*, 2023. doi:10.20944/preprints202308.1478.v1

[40] H. M. Abdelghany and K. Hooker, "Effective data sampling techniques for machine learning OPC in full chip production," *Proceedings of SPIE - The International Society for Optical Engineering*, vol. 11611, 2021. doi:10.1117/12.2586176

[41] C. Xie et al., "Effect of machine learning re-sampling techniques for imbalanced datasets in 18F-FDG PET-based radiomics model on prognostication performance in cohorts of head and neck cancer patients," *European Journal of Nuclear Medicine and Molecular Imaging*, vol. 47, no. 12, 2020. doi:10.1007/s00259-020-04756-4

[42] R. Gupta et al., "Diagnosis of Breast Cancer on Imbalanced Dataset Using Various Sampling Techniques and Machine Learning Models," *2021 International Conference on Digital Ecosystems and Technologies (DEST)*, 2021. doi:10.1109/DeSE54285.2021.9719398

[43] A. S. Dina et al., "Effect of Balancing Data Using Synthetic Data on the Performance of Machine Learning Classifiers for Intrusion Detection in Computer Networks," *IEEE Access*, vol. 10, 2022. doi:10.1109/ACCESS.2022.3205337

[44] H. Shi, Y. Xu, B. Ding, J. Zhou, and P. Zhang, "Long-Term Solar Power Time-Series Data Generation Method Based on Generative Adversarial Networks and Sunrise–Sunset Time Correction," *Sustainability*, vol. 15, no. 20, 2023. doi:10.3390/su152014920

[45] L. Mushunje, D. Allen, and S. Peiris, "Volatility and irregularity Capturing in stock price indices using time series Generative adversarial networks (TimeGAN)," *arXiv preprint arXiv:2311.12987*, 2023. doi:10.48550/arXiv.2311.12987

[46] C.-Y. Tai, W.-J. Wang, and Y.-M. Huang, "Using Time-Series Generative Adversarial Networks to Synthesize Sensing Data for Pest Incidence Forecasting on Sustainable Agriculture," *Sustainability*, vol. 15, no. 10, 2023. doi:10.3390/su15107834

[47] A. A. Purwita, A. Yesilkaya, and H. Haas, "Synthetic LiFi Channel Model Using Generative Adversarial Networks," *IEEE International Conference on Communications (ICC)*, 2022. doi:10.1109/ICC45855.2022.9838481

[48] S. Chattoraj, S. Pratiher, S. Pratiher, and H. Konik, "Improving Stability of Adversarial Li-ion Cell Usage Data Generation using Generative Latent Space Modelling," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021. doi:10.1109/ICASSP39728.2021.9413892

[49] C. Zhang and Y. Chen, "Time Series Anomaly Detection with Variational Autoencoders," 2019. doi:10.1109/ICMLA.2018.00207

[50] V. Fortuin, G. Rätsch, and S. Mandt, "Multivariate Time Series Imputation with Variational Autoencoders," 2019. doi:10.1109/ICMLA.2018.00207

[51] J. Li, W. Ren, and M. Han, "Mutual Information Variational Autoencoders and Its Application to Feature Extraction of Multivariate Time Series," 2022. doi:10.1142/s0218001422550059

[52] W. Todo, B. Laurent, J.-M. Loubes, and M. Selmani, "Dimension Reduction for time series with Variational AutoEncoders," 2022. doi:10.48550/arXiv.2204.11060

[53] G. G. González, P. Casas, A. Fernández, and G. Gómez, "Steps towards continual learning in multivariate time-series anomaly detection using variational autoencoders," 2022. doi:10.1145/3517745.3563033

[54] M. L. Garsdal, V. Sogaard, and S. M. Sørensen, "Generative time series models using Neural ODE in Variational Autoencoders," 2022. doi:Not available

[55] D. Li, D. Chen, L. Shi, B. Jin, J. Goh, and S.-K. Ng, "MAD-GAN: Multivariate Anomaly Detection for Time Series Data with Generative Adversarial Networks," 2019. doi:10.1007/978-3-030-30490-4_56

[56] W. Cheng, T. Ma, X. Wang, and G. Wang, "Anomaly Detection for Internet of Things Time Series Data Using Generative Adversarial Networks With Attention Mechanism in Smart Agriculture," 2022. doi:10.3389/fpls.2022.890563

[57] Z. Thompson, A. Downey, J. D. Bakos, and J. Wei, "Synthesizing Dynamic Time-series Data for Structures Under Shock Using Generative Adversarial Networks," 2022.

[58] K. Sarda, A. Yerudkar, and C. D. Vecchio, "Missing Data Imputation for Real Time-series Data in a Steel Industry using Generative Adversarial Networks," 2021. doi:10.1109/IECON48115.2021.9589716

[59] A. Takiddin, M. Ismail, U. Zafar, and E. Serpedin, "Deep Autoencoder-Based Anomaly Detection of Electricity Theft Cyberattacks in Smart Grids," 2022. doi:10.1109/JSYST.2021.3136683

[60] X. Jin, W. Gong, J. Kong, Y.-t. Bai, and T. Su, "PFVAE: A Planar Flow-Based Variational Auto-Encoder Prediction Model for Time Series Data," 2022. doi:10.3390/math10040610

[61] T. Kieu, B. Yang, C. Guo, R.-G. Cirstea, Y. Zhao, Y.-h. Song, and C. S. Jensen, "Anomaly Detection in Time Series with Robust Variational Quasi-Recurrent Autoencoders," 2022. doi:10.1109/icde53745.2022.00105

[62] A. Bouteska, M. Lavazza Seranto, p. hajek, and M. Z. Abedin, "Data-driven decadal climate forecasting using Wasserstein time-series generative adversarial networks," 2023. doi:10.1007/s10479-023-05722-7

[63] X. Hu, H. Zhang, D. Ma, and R. Wang, "Hierarchical Pressure Data Recovery for Pipeline Network via Generative Adversarial Networks," 2022. doi:10.1109/TASE.2021.3069003

[64] D. Wang, Y. Yan, R. Qiu, Y. Zhu, K. Guan, A. Margenot, and H. Tong, "Networked Time Series Imputation via Position-aware Graph Enhanced Variational Autoencoders," 2023. doi:10.1145/3580305.3599444

[65] F. Romanelli and F. Martinelli, "Synthetic Sensor Measurement Generation With Noise Learning and Multi-Modal Information," 2023. doi:10.1109/ACCESS.2023.3323038

[66] H. Qin, L. Su, C. Jiang, C. Zhang, G. Wu, and Y. Zhang, "Time Series Data Augmentation Algorithm Combining Deep Metric Learning and Variational Encoder," 2023. doi:10.1109/ICPSAsia58343.2023.10294708

[67] H. Li, S. Yu, and J. Príncipe, "Causal Recurrent Variational Autoencoder for Medical Time Series Generation," 2023. doi:10.48550/arXiv.2301.06574

[68] A. Siahkoohi, R. Morel, R. Balestriero, E. Allys, G. Sainton, T. Kawamura, and M. V. de Hoop, "Martian time-series unraveled: A multi-scale nested approach with factorial variational autoencoders," 2023. doi:10.48550/arXiv.2305.16189

[69] F. Altekrüger, P. Hagemann, and G. Steidl, "Conditional Generative Models are Provably Robust: Pointwise Guarantees for Bayesian Inverse Problems," 2023. doi:10.48550/arXiv.2303.15845

[70] P. Kashyap, C. Cheng, Y. Choi, and P. D. Franzon, "Generative Multi-Physics Models for System Power and Thermal Analysis Using Conditional Generative Adversarial Networks," 2023. doi:10.1109/EPEPS58208.2023.10314864

[71] X. Yuan, K. Chen, J. Zhang, W. Zhang, N. H. Yu, and Y. Zhang, "Pseudo Label-Guided Model Inversion Attack via Conditional Generative Adversarial Network," 2023. doi:10.48550/arXiv.2302.09814

[72] A. Elhagry, "Text-to-Metaverse: Towards a Digital Twin-Enabled Multimodal Conditional Generative Metaverse," 2023. doi:10.1145/3581783.3613432

[73] N. Vyas, S. Kakade, and B. Barak, "On Provable Copyright Protection for Generative Models," 2023. doi:10.48550/arXiv.2302.10870

[74] A. Heng and H. Soh, "Selective Amnesia: A Continual Learning Approach to Forgetting in Deep Generative Models," 2023. doi:10.48550/arXiv.2305.10120

[75] A. Tong, N. Malkin, G. Huguet, Y. Zhang, J. Rector-Brooks, K. Fatras, G. Wolf, and Y. Bengio, "Improving and generalizing flow-based generative models with minibatch optimal transport," 2023. doi:10.48550/arXiv.2302.00482

[76] D. Ye, X. Wang, and X. Chen, "Lightweight Generative Joint Source-Channel Coding for Semantic Image Transmission with Compressed Conditional GANs," 2023. doi:10.1109/ICCCWorkshops57813.2023.10233814

[77] A. A. Xu, S. Han, X. Ju, and H. Wang, "Generative Machine Learning for Detector Response Modeling with a Conditional Normalizing Flow," 2023. doi:10.48550/arXiv.2303.10148

[78] S. Liao, H. Ni, M. Sabaté-Vidales, L. Szpruch, M. Wiese, and B. Xiao, "Sig-Wasserstein GANs for conditional time series generation," 2023. doi:10.1111/mafi.12423

[79] H. Zhang, Z. Pang, J. Wang, and T. Li, "Few-shot Learning using Data Augmentation and Time-Frequency Transformation for Time Series Classification," 2023. doi:10.48550/arXiv.2311.03194

[80] M. F. Sikder, R. Ramachandranpillai, and F. Heintz, "TransFusion: Generating Long, High Fidelity Time Series using Diffusion Models with Transformers," 2023. doi:10.48550/arXiv.2307.12667

[81] Chen et al., "Data Augmentation for Pseudo-Time Series Using Generative Adversarial Networks," 2023. [Online]. Available: https://dblp.org/rec/conf/itat/SalmiJ23

[82] S. Crepey et al., "Anomaly Detection in Financial Time Series by Principal Component Analysis and Neural Networks," *Algorithms*, vol. 15, no. 10, p. 335, 2022. [Online]. Available: https://www.mdpi.com/1999-4893/15/10/335

[83] Z. Yang, Y. Li, G. Zhou, "TS-GAN: Time-series GAN for Sensor-based Health Data Augmentation," 2023. doi:10.1145/3583593

[84] P. Guo, H. Yang, A. Sano, "Empirical Study of Mix-based Data Augmentation Methods in Physiological Time Series Data," 2023. doi:10.1109/ICHI57859.2023.00037

[85] R. Peres, A. Gomes, J. Mendes, and S. Bento, "Simulation-Based Data Augmentation for the Quality Inspection of Structural Adhesive With Deep Learning," *IEEE Access*, vol. 9, pp. 44326-44335, 2021. doi:10.1109/ACCESS.2021.3069452

[86] H. Khosravi, S. Farhadpour, M. Grandhi, A. S. Raihan, S. Das, and I. Ahmed, "Strategic Data Augmentation with CTGAN for Smart Manufacturing: Enhancing Machine Learning Predictions of Paper Breaks in Pulp-and-Paper Production," *CoRR*, vol. abs/2311.09333, 2023. doi:10.48550/ARXIV.2311.09333

[87] J. Chung, B. Jang, "Accurate prediction of electricity consumption using a hybrid CNN-LSTM model based on multivariable data," 2022. doi:10.1371/journal.pone.0278071

[88] E. Branikas, P. Murray, G. West, "A Novel Data Augmentation Method for Improved Visual Crack Detection Using Generative Adversarial Networks," 2023. doi:10.1109/ACCESS.2023.3251988

[89] H. Qin, L. Su, C. Jiang, C. Zhang, G. Wu, and Y. Zhang, "Time Series Data Augmentation Algorithm Combining Deep Metric Learning and Variational Encoder," *Proc. ICPSAsia*, 2023. doi:10.1109/ICPSAsia58343.2023.10294708

[90] Z. Cai, W. Ma, X. Wang, H. Wang, and Z. Feng, "The Performance Analysis of Time Series Data Augmentation Technology for Small Sample Communication Device Recognition," *IEEE Transactions on Robotics*, vol. 39, no. 1, pp. 5-15, 2023. doi:10.1109/TR.2022.3178707

[91] B. Shen, L. Yao, X. Jiang, Z. Yang, and J.-s. Zeng, "Time Series Data Augmentation Classifier for Industrial Process Imbalanced Fault Diagnosis," *Proc. DDCLS*, 2023. doi:10.1109/DDCLS58216.2023.10166336

[92] Y. Gao, C. A. Ellis, V. Calhoun, and R. L. Miller, "Improving age prediction: Utilizing LSTM-based dynamic forecasting for data augmentation in multivariate time series analysis," *arXiv preprint arXiv:2312.08383*, 2023.

[93] A. Wilf, A. T. Xu, P. Liang, A. Obolenskiy, D. Fried, and L.-P. Morency, "Comparative Knowledge Distillation," *arXiv preprint arXiv:2311.02253*, 2023.

[94] K. Rath, D. Rügamer, B. Bischl, U. von Toussaint, C. Rea, A. D. Maris, R. Granetz, and C. Albert, "Data augmentation for disruption prediction via robust surrogate models," *Journal of Plasma Physics*, vol. 88, no. 4, 2022. doi:10.1017/S0022377822000769

[95] Y. Kwak and J.-g. Huh, "Random Augmentation Technique for Mitigating Overfitting in Neural Networks for Financial Time Series Forecasting," *Journal of Korean Data Analysis Society*, vol. 25, no. 5, pp. 1653-1664, 2023. doi:10.37727/jkdas.2023.25.5.1653

[96] M. Li and B. C. Lovell, "End to End Generative Meta Curriculum Learning for Medical Data Augmentation," in *Proc. ICIP*, 2022. doi:10.1109/ICIP49359.2023.10222093

[97] J. C. Lin and F. Yang, "Data Augmentation for Industrial Multivariate Time Series via a Spatial and Frequency Domain Knowledge GAN," in *Proc. AdCONIP*, 2022. doi:10.1109/AdCONIP55568.2022.9894177

[98] J. Yuan, R. Tang, X. Jiang, and X. Hu, "LLM for Patient-Trial Matching: Privacy-Aware Data Augmentation Towards Better Performance and Generalizability," *arXiv preprint arXiv:2303.16756*, 2023.

[99] M. Rahul and S. Chiddarwar, "A causality-inspired data augmentation approach to cross-domain burr detection using randomly weighted shallow networks," *International Journal of Machine Learning and Cybernetics*, vol. 14, no. 2, pp. 569-582, 2023. doi:10.1007/s13042-023-01891-w

[100] A. Gong, X. Zhang, Y. Wang, Y. Zhang, and M. Li, "Hybrid Data Augmentation and Dual-Stream Spatiotemporal Fusion Neural Network for Automatic Modulation Classification in Drone Communications," *Drones*, vol. 7, no. 6, 2023. doi:10.3390/drones7060346

[101] T. Döhmen, M. Hulsebos, C. Beecks, and S. Schelter, "GitSchemas: A Dataset for Automating Relational Data Preparation Tasks," in *Proc. ICDEW*, 2022. doi:10.1109/icdew55742.2022.00016

[102] H. Chen and Y. Ji, "Improving the Interpretability of Neural Sentiment Classifiers via Data Augmentation," in *Proc. EMNLP-IJCNLP*, 2019.

[103] P. Katiyar and A. Khoreva, "Improving Augmentation and Evaluation Schemes for Semantic Image Synthesis," *arXiv preprint arXiv:2011.12636*, 2020.

[104] F. Xie, H. Wen, J. Wu, W. Hou, H.-h. Song, T. Zhang, R. Liao, and Y. Jiang, "Data Augmentation for Radio Frequency Fingerprinting via Pseudo-Random Integration," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 4, no. 5, pp. 594-604, 2020. doi:10.1109/TETCI.2019.2907740

# Utilizing UAV Data for Neural Network-based Classification of Melon Leaf Diseases in Smart Agriculture

Siti Nur Aisyah Mohd Robi[1], Norulhusna Ahmad[2], Mohd Azri Mohd Izhar[3], Hazilah Mad Kaidi[4],
Norliza Mohd Noor[5],
Razak Faculty of Technology and Informatics, Universiti Teknologi Malaysia (UTM) Kuala Lumpur, Malaysia[1,2,3,4]
ATES Sdn. Bhd., Kuala Lumpur[5]

*Abstract*—**Integrating unmanned aerial vehicle (UAV) technology with plant disease detection is a significant advancement in agricultural surveillance, marking the beginning of a transformational era characterised by innovation. Traditionally, farmers have had to rely on manual visual inspections to identify melon leaf diseases, which proves to be a time-consuming and costly process in terms of labour. This paper aims to use UAV technology for plant disease detection to achieve notable progress in agricultural surveillance. Incorporating UAV technology, specifically utilising the You Only Look Once version 8 (YOLOv8) deep-learning model, is revolutionary in precision agriculture. This study uses UAV imagery in precision agriculture to explore the utility of YOLOv8, a powerful deep-learning model, for detecting diseases in melon leaves. The labelled dataset is created by annotating disease-affected areas using bounding boxes. The YOLOv8 model has been trained using a labelled dataset to detect and classify various diseases accurately. Following the training, the performance of YOLOv8 stands out significantly compared to other models, boasting an impressive accuracy of 83.2%. This high level of accuracy underscores its effectiveness in object detection tasks and positions it as a robust choice in computer vision applications. It has been shown that rigorous evaluation can help find diseases, which suggests that it could be used for early intervention in precision farming and to change how crop management systems work. This has the potential to assist farmers in promptly identifying and addressing plant issues, hence altering their crop management practices.**

*Keywords*—*Smart agriculture; plant disease; melon leaf disease; image processing; neural network; UAV*

## I. INTRODUCTION

The melon, scientifically known as Cucumis melo L., is a significant horticultural commodity thriving in Malaysia, categorized within the Cucurbitaceae family. Its cultivation is widespread globally, particularly on subtropical and tropical regions [1], [2]. Melons are esteemed for their delightful sweet tastes, crisp textures, and distinct fragrances [3]. The flesh of the melon is also very good for you because it is full of ascorbic acid (vitamin C), which is a water-soluble vitamin that is known to be one of the safest and most effective nutrients [4], [5]. Numerous farmers extensively cultivate melon, a high-value fruit commodity. Cultivating melons is challenging due to the prevalence of various diseases associated with melon plants. Leaf diseases in melon plants result in economic losses for melon growers. Melon plant diseases are classified into two categories according to their causes: insects and viruses.

One of the insects is called a leaf miner. Various polyphagous leafminer flies pose a potential threat to vegetable crops, and occasionally, even melons are susceptible to these pests. Then, mines show up on the leaflets. The worst-affected leaves may become yellow, wilt, and dry out. These leaves might occasionally contain up to 20 larvae per. Thus, during an infestation, a plant's photosynthetic activity, growth, and yields can all be significantly decreased [6]. Other than leaf miners, aphids are also one of the insects that can form colonies on the young leaflets of melon leaf. They usually establish colonies when they develop on melon. They are particularly dangerous since they can spread many viruses. Nutritional punctures cause chlorotic punctures, which can deform young, rolled-up, and somewhat bloated leaves.

To tackle this challenge, image processing, machine learning (ML), and deep learning (DL) methods offer a solution for categorising plant disease levels on melon leaves, aiding farmers in effectively managing these issues. These techniques have been widely employed in identifying, detecting, and classifying different types of leaf diseases. Scholars have initiated investigations into applying deep learning models for plant detection and counting. This includes the utilization of popular models such as you only look once (YOLO) [7], faster region-based convolutional neural network (Faster R-CNN) [8], and EfficientDet [9]. Several scholars have also implemented a series of enhancements to achieve the objectives of plant detection and counting jobs [10], [11]. Different methods were used to classify leaf diseases: DenseNet and Inception categorized four diseases for bananas, with DenseNet showing better accuracy at 84.73% [12]. Grape leaves were classified into healthy and leaf spot categories using deep forest, achieving 96.25% accuracy [13]. Cucumber leaf diseases were segmented to identify disease points, reaching 97.23% accuracy using an improved saliency method and deep feature selection [14]. Detecting and categorising leaf diseases involves extracting features, which are then used for classification [15].

This research presents a novel approach to disease detection in melon plants through a neural network using drone imagery and an effective method known as YOLO. It is a novel strategy for handling melon problems from above, assisting farmers in more accurately identifying and controlling crop diseases. The contributions of this paper are as follows:

- The research proposes a unique methodology for disease detection in melon plants by employing a neural network trained on drone imagery. This inno-

vative approach aims to address the identified gap in the literature and contribute to advancing agricultural monitoring.

- The study employs the YOLOv8 and YOLOv5 methods, demonstrating their effectiveness as an efficient and accurate tool for identifying diseases in melon plants.

- The research contribution lies in its potential to significantly improve the overall management of melon plants by introducing a novel combination of unmanned aerial vehicle (UAV) imagery and the YOLO method.

As a result, this paper presents a novel method for identifying diseases in melon plants, addressing the gap in current investigations using the YOLO model. This method can significantly improve the process by which farmers detect and manage infections in melon crops, representing a notable advancement in agricultural techniques.

## II. Related Works

In precision agriculture, drones have been applied in various ways, and new applications are always being investigated. Numerous drone applications have been created for various uses, including soil analysis, pest detection, crop yield estimation, yield spraying, water stress detection, land mapping, plant nutrient deficiency identification, livestock control, weed detection, and protection of agricultural products [16]. Using UAVs to detect plant leaf diseases has grown in popularity over the past few years due to the industry's rapid growth in machine vision and UAV manufacture [17].

The effective use of DL technology in plant disease categorization in recent years has given researchers a fresh perspective on the topic. Traditionally, disease diagnosis in farming has depended on unaided eye observation, which is costly, time-consuming, and highly skilled [18]. It is possible that deep learning methods could help solve problems in feature extraction, classification, and expert system development. This could help farmers grow better fruit plants that produce more fruit. Models like DenseNet-121 [19], ResNet-50 [20], and MobileNet [21] are well-known and have been used in many previous studies to find and classify images in the field of diagnosing and identifying plant diseases. Sladojevic et al. introduced a method for identifying plant diseases utilizing a Convolutional Neural Network (CNN) within the Caffe DL framework [22]. They gathered images from diverse origins and employed data augmentation methods such as affine transformation, perspective transformation, and rotation to create additional images.

The YOLO model has gained considerable attention due to its remarkable combination of accuracy and speed. Regression-based object detection models commonly used are Single Shot Multi-box Detector (SSD) [23], and YOLO [24]. YOLO is a basic neural network that can simultaneously predict bounding box coordinates and related class probabilities. Also, YOLO frame detection is seen as a regression problem because it finds targets from start to finish without the need for a complicated pipeline [24], making it very efficient. Moreover, YOLO outperforms other real-time systems regarding mean

average precision (mAP) [25]. Goyal et al. proposed a model based on the YOLOv5 object detection system to sort fruit for fruit detection and quality detection [26]. For fruit detection, the model's mAP was 92.80% in the first stage and 95.60% and 93.10% for apples and bananas, respectively, in the second stage. For pear counters, Parico and Ahamed employed depth sorting and the YOLOv4 model to recognize and count pear fruit in real time [11]. YOLOv8, the most recent iteration in the YOLO series, not only retains its predecessors' strengths but surpasses them, thereby emerging as a powerful instrument for professionals in plant science.

Besides, no existing methods are designed to detect disease in this melon, representing a research gap in using DL with UAV images for melon diseases. Although DL methods and UAV imagery have been utilized in research to detect diseases in other plants, melon diseases have not received as much attention as they should. The realised gap highlights the lack of thorough investigation into the potential advantages of using UAV images for disease detection in melon crops. This highlights the need for targeted research in this specific area. Utilising deep learning to identify plant diseases may overcome the drawbacks of manually selecting disease spot characteristics. This approach enhances the objectivity of plant disease feature extraction and accelerates research efficiency and technological transition.

## III. Proposed Approach

This study used a neural network model and an image processing technique to develop a method for identifying melon leaf diseases. Fig. 1The proposed approach. shows the proposed approach. The details of the proposed method will be discussed in the next section.

### A. Data Acquisition

A dataset of melon leaf images was gathered from KMK Agro Global Sdn. Bhd., Banting, Selangor, with the leaves seen in their natural environment under a controlled greenhouse. Moreover, the dataset has been extracted specifically for analysing colour. The flowchart in Fig. 2The flowchart of the system for disease detection using UAV images. describes a systematic process for combining UAV-captured imagery and the DL model YOLOv8 to detect diseases in melon leaves. First, information is gathered using a UAV to take recordings of the melon greenhouse. Pre-processing operations are performed on the collected data, such as consistent image expansion, image removal, and extracting video recordings into their component frames for analysis. The images are then labelled by highlighting spots that indicate illnesses on the rock melon leaves and annotating locations of interest with bounding boxes. The labelled dataset becomes the foundation for training the YOLOv8 model. During this phase, the model learns to recognise and classify different diseases affecting the leaves. After training, the model is rigorously tested using a different set of images. This evaluation step uses performance metrics like accuracy, precision, and recall to see how well the model can reliably find and classify illnesses in rock melon leaves.

The melon dataset was collected using a high-resolution UAV, DJI Mavic Air 2s. Fig. 3 illustrates the UAV, DJI Mavic

Fig. 1. The proposed approach.



Fig. 2. The flowchart of the system for disease detection using UAV images.



Fig. 3. DJI Mavic Air 2s.



Fig. 4. The layout and drone flying direction in the KMK Agro Global Sdn. Bhd. greenhouse for data collection.

Air 2s used during data collection. The UAV specifications are shown in Table IDrone Specifications. When engaging in the photographic documentation of melon leaves, it becomes imperative to factor in technical intricacies. This involves maintaining a precise distance centred on the leaf object, within 15 cm to 20 cm. Ensuring that the leaf object remains well-contained within the camera frame is crucial. From the aerial perspective, the drone's movement will be orchestrated upward and downward along the plant, meticulously scrutinising for any signs of disease. The height of the plant varied from 1.5 m to 2.0 m. The drone captured images of healthy plants and plants with diseases. The images captured must be within 20 cm of the plant so that the leaves are visible in the drone's field of view (FOV).

Fig. 4 illustrates the greenhouse layout and the direction in which the drone is flying. The recording is in 4K and at normal speed to ensure the high quality of the images. The UAV flies facing the plant and moves along the row. Throughout the data collection, a UAV captured data in 4K resolution during a 30-minute recording session. The video will be broken down into frames to extract images of the melon plant at five-second intervals. Fig. 5 depicts some of the images captured by UAV.

### B. Data Augmentation

After data collection, frames were extracted from the video every five seconds to obtain appropriate images for training. The dataset collected from the farm for melon plant disease exhibits an imbalance, which may compromise the

TABLE I. DRONE SPECIFICATIONS

| No | Feature | Specification |
|----|---------|---------------|
| 1. | Drone | DJI Mavic Air 2s |
| 2. | Video Resolution | 4K: 3840 x 2160 at 24/25/30/48/50/60fps |
| 3. | Max Flight Time | Up to 34 Minutes |
| 4. | Camera Sensor | 1/2-inch CMOS, 48 MP |



Fig. 5. Images of melon plants collected from UAVs before dataset training.

accuracy of the YOLO model. Data augmentation enhances the model's performance by generating diverse variations of the training data. This reduces the problem of overfitting and enhances the model's capacity to form generalisations. In this study, ImageDataGenerator by Keras library is used for data augmentation. Using the ImageDataGenerator class in Keras makes it easy to set up and apply random transformations to image data, such as rotations, shifts, flips, and normalisation operations. These changes can be made without interrupting the training pipeline, which makes the model more flexible and good at what it does.

### C. Data Annotation

After the data augmentation, the preprocessing stage aimed to identify plants with diseases for the training process. Once all the images were selected, the labelling process began. When labelling images in computer vision, bounding boxes annotate objects or regions of interest by enclosing them with rectangles or other shapes. Neural network models find it easier to locate, localise, and recognise items when using these bounding boxes, which accurately show the location and bounds of certain objects. Labelling is challenging as it involves addressing imbalances in disease images, which could impact training accuracy. Balancing diseased plant images with normal ones ensures the YOLO model functions effectively. Fig. 6Data labelling process: Annotated markings highlighting disease-affected areas on melon leaf. shows the labelling process of the melon leaf.



Fig. 6. Data labelling process: Annotated markings highlighting disease-affected areas on melon leaf.

For this melon plant, five classes were utilized in the labelling process. These classes include normal, unknown, mosaic, leafminer, and aphid. All these types of diseases commonly affect melon plants. Fig. 7(a) Aphid (b) Leafminer (c) Mosaic (d) Unknown. displays the diseases that typically affect melon plants. For a normal melon plant, the leaves are green and devoid of white or yellow spots.



Fig. 7. (a) Aphid (b) Leafminer (c) Mosaic (d) Unknown.

### D. Training the Dataset

The dataset has been divided into three distinct sections: training, testing, and evaluation, each accounting for 80%,

TABLE II. Parameters before Dataset Training

| Training Images | 1200 |
|---|---|
| Test Images | 150 |
| Size Images | 640x640 |
| Epoch | 100 |
| Class | 5 classes |
| GPU | NVIDIA GPU |

10%, and 10%, respectively. Before the training phase, the images undergo augmentation techniques to increase their quantity. This involves resizing them uniformly to 640x640 and introducing rotations within the range of +15° and -15°. The training dataset comprises 1200 images, while 150 images have been allocated for testing and evaluation. This process sets the stage for the model to excel in identifying elusive diseases nestled within melon leaves.

*1) YOLOv8 Model:* YOLOv8 is the latest cutting-edge model within the YOLO series, suitable for object detection, image classification, and instance segmentation tasks. The influential and industry-shaping YOLOv5 model's creators, Ultralytics, are also responsible for creating the YOLOv8, a significant advancement in this field. Consider utilizing YOLOv8 for your upcoming computer vision endeavour for several compelling reasons. Firstly, its accuracy, assessed via common objects in context (COCO) and Roboflow 100 metrics, stands notably high. Secondly, YOLOv8 boasts a range of developer-friendly features, including an intuitive CLI and a well-structured Python package, enhancing usability. A robust community within the YOLO framework, particularly around the YOLOv8 version, supports the model. This means that people who work in computer vision can get much help and advice. Notably, YOLOv8 demonstrates robust performance on COCO benchmarks, exemplified by the YOLOv8m model achieving a 50.2% mAP. Table IIParameters before Dataset Training displays the comprehensive set of parameters employed specifically for training the YOLOv8 model.

## IV. Result and Discussion

After the training process, the model is tested and evaluated. Performance metrics considered in this study are mean average precision (mAP), precision, and recall. The results are shown in Table IIIPerformance Evaluation for Plant Disease. From Table IIIPerformance Evaluation for Plant Disease, YOLOv8 outperforms YOLOv5. The dataset used during the training is the same. The superiority of YOLOv8 over YOLOv5 is evident through substantial enhancements. YOLOv8 exhibits an mAP of 83.2%, precision at 84.3%, and a recall of 73%. YOLOv8 performs better than YOLOv5 because of several significant improvements and optimizations. These enhancements could include improved network architectures, feature extraction strategies, sophisticated training approaches, or hyperparameter adjustments. It is possible that YOLOv8's more complex or effective backbone architecture allowed it to extract more significant features from the data.

Furthermore, YOLOv8 experienced extensive training procedures using bigger and more varied datasets, which improved its capacity to generalize and precisely identify objects and

resulted in greater precision and recall scores. Compared to YOLOv5, YOLOv8 performs better overall because of these enhancements combined. Once the model is ready, it undergoes rigorous testing to evaluate its performance and accuracy. Fig. 8Results of using YOLOv8 on test images to spot diseases in melons. depicts the output of the detection process, showcasing the model's performance. The evaluation of the system for categorization included metrics such as mAP, precision, and recall to measure its efficacy. MAP is a metric used to evaluate the performance of object detection models. It measures the average precision of an algorithm across multiple classes or object categories. It considers the precision and recall of the model's predictions, offering a comprehensive assessment of how accurately and completely the model detects objects in an image across different categories.

After the training, other parameters can be employed to evaluate the effectiveness of YOLOv8 detection algorithms. Once the dataset training is over, there is a difference in the accuracy of class identification. In particular, the classes related to mosaic illness and the unknown condition in this investigation showed noticeably lower accuracy scores of are 75% and 76%, respectively, as shown in Fig. 9The outcomes observed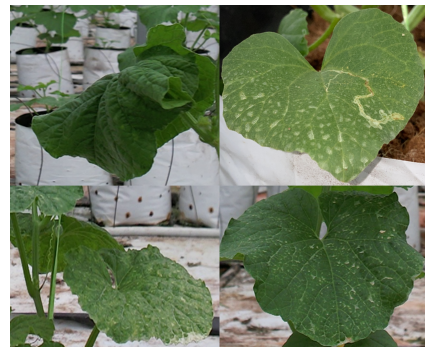 for each class after the training phase.. This decreased accuracy is because several diseases have remarkably similar traits, making it difficult for the model to discriminate between them. These particular diseases are difficult to classify accurately due to their intricate visual traits and similarities, which is why these classes' accuracy levels are lower. The area under the precision-recall curve at different detection thresholds is called AP. The mAP shows how accurate the system is for each of the $n$ object classes.



Fig. 8. Results of using YOLOv8 on test images to spot diseases in melons.



Fig. 9. The outcomes observed for each class after the training phase.

The mAP% can be calculated using the equation [27]

TABLE III. PERFORMANCE EVALUATION FOR PLANT DISEASE

| Model | Number of Images | mAP | Precision | Recall |
|-------|-----------------|------|-----------|--------|
| YOLOv5 | 1200 | 72.7% | 83.3% | 65.7% |
| YOLOv8 | 1200 | 83.2% | 84.3% | 73% |

$$\text{mAP} = \frac{\sum_i^n AP_i}{n}. \tag{1}$$

$$\text{AP} = \int_0^1 x\,dy \tag{2}$$

Precision and recall can be measured using true positive (TP), true negative (TN), false positive (FP), and false-negative (FN) indicators. The equation to calculate precision $x$, and recall $y$ are given by:

$$x = \frac{\text{TP}}{\text{TP} + \text{FP}}, \tag{3}$$

and

$$y = \frac{\text{TP}}{\text{TP} + \text{FN}}, \tag{4}$$

where $x$ is precision and $y$ is recall.

Other than YOLOv8, this study compares its performance with YOLOv5. The YOLOv8 proposed method leaves the use of predefined anchor boxes and instead employs an anchor-free strategy to achieve enhanced item localization accuracy, particularly for smaller objects. The Path Aggregation Network (PANet) integrates several network-level features, enhancing detection accuracy using multi-scale contextual information. Fig. 10YOLOv5 training development graph, highlighting recall, precision, and accuracy metrics during training. and Fig. 11YOLOv8 training development graph, highlighting recall, precision, and accuracy metrics during training. show the training graph for YOLO. The training process involves iterating through the dataset 100 times, each iteration known as an epoch. During these 100 epochs, the model learns and refines its understanding of the data, gradually improving its performance and accuracy through repeated exposure to the information provided in the dataset.

Moreover, the convolutional block attention module (CBAM) has been enhanced to improve the feature extraction process. The dynamic adjustment of feature importance achieves this by effectively suppressing noise and improving the clarity of distinctions. The efficient backbone network of YOLOv8 successfully preserves accuracy by reducing parameters and enhancing inference performance. The proposed approach effectively separates the responsibilities of object prediction and categorization, resulting in improved precision. The system attains accelerated convergence and enhanced stability by employing network pruning, varied data augmentation techniques, mixed precision training, and an enhanced training framework. The unified architecture of YOLOv8 enables its compatibility with many vision tasks, establishing it as a



Fig. 10. YOLOv5 training development graph, highlighting recall, precision, and accuracy metrics during training.



Fig. 11. YOLOv8 training development graph, highlighting recall, precision, and accuracy metrics during training.

robust and versatile tool for applications involving object identification and image recognition.

## V. LIMITATIONS OF THIS STUDY

The study highlights certain limitations that should be recognised. Firstly, a limited annotated dataset was utilised for training the YOLOv8 model. Using a restricted dataset raises concerns regarding the possibility of overfitting, wherein the model may exhibit good performance on the training data but encounter difficulties in efficiently applying its knowledge to new, real-world farming situations. Furthermore, in this study, the influence of environmental factors was considered by conducting data collection in a controlled atmosphere. Climatic conditions might impact the quality of UAV footage and the precision of the YOLOv8 model's forecasts, underscoring the necessity to tackle these obstacles for real-world implementations. Furthermore, it is necessary to carefully examine the possible impact of environmental variables, such as changes in lighting, on the performance of the UAV and YOLOv8 model.

## VI. CONCLUSION

In summary, this study has used UAV footage to highlight the strong performance of YOLOv8 in detecting diseases in rock melon leaves. The deep learning model exhibits notable levels of accuracy and efficiency, highlighting its potential as a useful asset in precision agriculture. While acknowledging problems like limited datasets and the effect of changes in the environment on the performance of models, the study has

highlighted the strong features of YOLOv8 as a major step forward in finding illnesses in agricultural settings. Overall, YOLOv8 emerges as a pivotal technological advancement, promising significant enhancements and advancements in agricultural practices.

This study offers numerous tangible benefits for the agricultural industry. Combining UAV data with a neural network-based classification system greatly improves the identification of melon leaf illnesses, allowing farmers to detect problems at their initial stages. The efficiency of this technology is especially advantageous for monitoring extensive agricultural regions, resulting in time and labour savings compared to conventional manual techniques. The neural network enables rapid identification, timely intervention, and effective disease management. The technology's capacity to scale allows it to be easily adjusted for large-scale farming operations, and the data-driven decision-making process provides farmers with vital knowledge to manage crops effectively. The research has the potential to fundamentally transform disease control, resulting in higher crop productivity, enhanced quality, and the adoption of more sustainable farming methods in smart agriculture.

Further studies may considerably improve the performance of the neural network-based classification system using other types of DL models. Actively focusing on diversity in the training dataset is one important approach. This study can incorporate samples from various geographic regions, meteorological conditions, and growing seasons to create a more comprehensive dataset. This diversity would strengthen the model's robustness and generalizability across many agricultural contexts, as well as its capacity to adjust to various environmental conditions. Then, further research could examine how UAV data can be integrated with other cutting-edge sensing technologies. For example, drones have been added to monitoring systems. This multidisciplinary strategy might result in a more comprehensive and precise disease detection solution for smart agriculture. Combining different sensing technologies could lead to a more comprehensive understanding of crop health and, ultimately, a more advanced and efficient precision agriculture system.

### ACKNOWLEDGMENT

### REFERENCES

[1] H. Kesh and P. Kaushik, "Advances in melon (Cucumis melo L.) breeding: An update," *Scientia Horticulturae*, vol. 282, p. 110045, 2021.

[2] P. Rolim, L. M. J. Seabra, and G. R. de Macêdo, "Melon by-products: Biopotential in human health and food processing," *Food Reviews International*, vol. 36, pp. 15 – 38, 2020.

[3] N. N. Nguyen, K. YongSham, P. JaeRyoung, and S. SungChur, "Development of a core set of ssr markers for cultivar identification and seed purity tests in oriental melon (cucumis melo l. var. makuwa)," *korean Journal of Horticultural Science&Technology*, 2019.

[4] M. Wen, S. Yang, L. Huo, P. He, X. Xu, C. Wang, Y. Zhang, and W. Zhou, "Estimating Nutrient Uptake Requirements for Melon Based on the QUEFTS Model," *Agronomy*, vol. 12, no. 1, 2022.

[5] E. Evana and M. Barek, "Determination of vitamin c (ascorbic acid) contents in two varieties of melon fruits (cucumis melo l.) by iodometric titration," *Fullerene Journal of Chemistry*, vol. 6, no. 2, pp. 143–147, 2021.

[6] S. Kumar, V. Bojan, S. V. Krishnamoorthy, L. Rajendran, B. Vinothkumar, S. S. Monica, J. K. Sathishkumar, and S. V. Krishnamoorthy, "A review on management of leafminer in horticultural crops," *Journal of Entomology and Zoology Studies*, vol. 9, no. 2, pp. 1204–1213, 2021.

[7] Z. F. Xu, R. S. Jia, Y. B. Liu, C. Y. Zhao, and H. M. Sun, "Fast method of detecting tomatoes in a complex scene for picking robots," *IEEE Access*, vol. 8, pp. 55 289–55 299, 2020.

[8] Y. Liu, C. Cen, Y. Che, R. Ke, Y. Ma, and Y. Ma, "Detection of Maize Tassels from UAV RGB Imagery with Faster R-CNN," *Remote Sensing*, vol. 12, no. 2, 2020. [Online]. Available: https://www.mdpi.com/2072-4292/12/2/338

[9] Y. Wang, Y. Qin, and J. Cui, "Occlusion Robust Wheat Ear Counting Algorithm Based on Deep Learning," *Frontiers in Plant Science*, vol. 12, 2021. [Online]. Available: https://www.frontiersin.org/articles/10.3389/fpls.2021.645899

[10] D. Lu, J. Ye, Y. Wang, and Z. Yu, "Plant Detection and Counting: Enhancing Precision Agriculture in UAV and General Scenes," *IEEE Access*, vol. 11, no. September, pp. 116 196–116 205, 2023.

[11] A. I. Parico and T. Ahamed, "Real Time Pear Fruit Detection and Counting Using YOLOv4 Models and Deep SORT," 2021.

[12] A. Ridhovan, A. Suharso, and C. Rozikin, "Disease detection in banana leaf plants using densenet and inception method," *Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi)*, 2022. [Online]. Available: https://api.semanticscholar.org/CorpusID:252675478

[13] J. and Arora, U. Agrawal, and P. Sharma, "Classification of Maize leaf diseases from healthy leaves using Deep Forest," *Journal of Artificial Intelligence and Systems*, vol. 2, no. 1, pp. 14–26, 2020.

[14] M. A. Khan, T. Akram, M. Sharif, K. Javed, M. Raza, and T. Saba, "An automated system for cucumber leaf diseased spot detection and classification using improved saliency method and deep features selection," *Multimedia Tools and Applications*, vol. 79, no. 25, pp. 18 627–18 656, 2020. [Online]. Available: https://doi.org/10.1007/s11042-020-08726-8

[15] L. C. Ngugi, M. Abelwahab, and M. Abo-Zahhad, "Recent advances in image processing techniques for automated leaf pest and disease recognition – A review," *Information Processing in Agriculture*, vol. 8, no. 1, pp. 27–51, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2214317320300196

[16] P. Zambrano, F. Calderon, H. Villegas, J. Paillacho, D. Pazmiño, and M. Realpe, "Uav remote sensing applications and current trends in crop monitoring and diagnostics: A systematic literature review," in *2023 IEEE 13th International Conference on Pattern Recognition Systems (ICPRS)*, 2023, pp. 1–9.

[17] W. Li, X. Yu, C. Chen, and Q. Gong, "Identification and localization of grape diseased leaf images captured by UAV based on CNN," *Computers and Electronics in Agriculture*, vol. 214, no. September, p. 108277, 2023.

[18] G. K. Sandhu and R. Kaur, "Plant Disease Detection Techniques: A Review," *2019 International Conference on Automation, Computational and Technology Management, ICACTM 2019*, pp. 34–38, 2019.

[19] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," 2018.

[20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2015.

[21] E. Elfatimi, R. Eryigit, and L. Elfatimi, "Beans leaf diseases classification using mobilenet models," *IEEE Access*, vol. 10, pp. 9471–9482, 2022.

[22] S. Sladojevic, M. Arsenovic, A. Anderla, D. Culibrk, and D. Stefanovic, "Deep Neural Networks Based Recognition of Plant Diseases by Leaf Image Classification," *Computational Intelligence and Neuroscience*, vol. 2016, p. 3289801, 2016.

[23] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, *SSD: Single Shot MultiBox Detector*. Springer International Publishing, 2016, p. 21–37.

[24] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2016-Decem, pp. 779–788, 2016.

[25] C.-K. Chang, C. Siagian, and L. Itti, "Mobile robot vision navigation and localization using gist and saliency," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2010, pp. 4147–4154.

[26] K. Goyal, P. Kumar, and K. Verma, "AI-based fruit identification and quality detection system," *Multimedia Tools and Applications*, vol. 82, no. 16, pp. 24 573–24 604, 2023.

[27] N. Zainab, H. Afzal, T. Al-shehari, M. Al-razgan, M. Zakria, M. J. Hyder, and R. Nawaz, "Detection and Classification of Temporal Changes for Citrus Canker Growth Rate using Deep Learning," *IEEE Access*, vol. PP, p. 1, 2023.

# Guiding 3D Digital Content Generation with Pre-Trained Diffusion Models

Jing Li[1¶], Zhengping Li[2¶*], Peizhe Jiang[3], Lijun Wang[4], Xiaoxue Li[5*], Yuwen Hao[6]

School of Information, North China University of Technology, Beijing, China 100144[1, 2, 4]

Shenzhen Renrenzhuang Technology Co., Ltd, Shenzhen, China 518000[3]

Beijing Key Laboratory of Disaster Rescue Medicine,

Medical Innovation Research Division of the Chinese PLA General Hospital, Beijing, China 100853[5, 6]

[¶]The authors have an equal contribution

*Abstract*—The production technology of 3D digital content involves multiple stages, including 3D modeling, simulation animation, visualization rendering, and perceptual interaction. It is not only the core technology supporting the creation of 3D digital content but also a key element in enhancing immersive application experiences in virtual reality and the metaverse. A primary focus in computer vision and computer graphics research has been on how to create 3D digital content that is efficient, convenient, controllable, and editable. Currently, producing high-quality 3D digital content still requires significant time and effort from a large number of designers. To address this challenge, leveraging artificial intelligence-generated methods to break down production barriers has emerged as an effective strategy. With the substantial breakthroughs achieved by diffusion models in the field of image generation, they also demonstrate tremendous potential in 3D digital content generation, potentially becoming a foundational model in this area. Recent studies have shown that diffusion model-based techniques for generating 3D digital content can significantly reduce production costs and enhance efficiency. Therefore, it is essential to summarize and categorize existing methods to facilitate further research. This paper systematically reviews 3D digital content generation methods, introducing related 3D representation techniques and focusing on 3D digital content generation schemes, algorithms, and pipeline based on diffusion models. We perform a horizontal comparison of different approaches in terms of generation speed and quality, deeply analyze existing challenges, and propose viable solutions. Furthermore, we thoroughly explore future research themes and directions in this domain, aiming to provide guidance and reference for subsequent research endeavors.

*Keywords*—*3D Digital content; computer vision; artificial intelligence; diffusion models; 3D representation*

## I. INTRODUCTION

Humans describe the world through text, comprehend it through images, and experience and interact with it in a three-dimensional (3D) format. Therefore, generative models have found widespread application in numerous aspects of life, playing a significant role in advancing human society. Research in recent years has mainly focused on text generation [1], [2], [3], [4] and image generation [5], [6], [7], [8]. Text generation is typically used for language tasks such as translation and question-answering, while image generation often involves creating visuals based on textual prompts. The generation of 3D digital content has not yet achieved the extraordinary capabilities seen in the domains of text and image generation. Therefore, there is still a need to continue to promote related research on 3D digital content generation.

3D digital content is extensively utilized in fields such as film, architecture, virtual and augmented reality. However, the current mainstream production of 3D digital content relies heavily on 3D designers, leading to remarkably low production efficiency and high entry barriers. Consequently, employing artificial intelligence (AI) to generate 3D digital content can significantly enhance production efficiency, reduce industry barriers, and foster the development of related fields.

Zero-shot image models [9] are trained using hundreds of millions of graphics data, which is difficult to achieve in the 3D domain. Table I presents a comparison between the data volumes of mainstream 3D and 2D datasets. Conventional 3D digital content generation methodologies predominantly utilize 3D datasets for training specific generative models [10], [11]. The advantage of this method lies in its ability to generate 3D objects with consistent geometry. However, it is limited by the current lack of sufficiently large 3D datasets and the absence of efficient 3D digital content generation architectures, as well as the computational power needed for their training. Therefore, it is difficult for this 3D digital content generation method to achieve a breakthrough in the short term. In light of this, this paper focuses on using pre-trained diffusion models [7], [12], [13] to supervise the generation of 3D digital content.

TABLE I. COMPARISON OF 3D DATASETS AND 2D DATASETS

| 3D | | | 2D | |
|---|---|---|---|---|
| Dataset | Full Mesh | Objects | Dataset | Images |
| ShapeNet [14] | ✓ | 51K | ImageNet [15] | 14M |
| AKB-48 [16] | ✓ | 2K | COCO [17] | 330K |
| OmniObject3D [18] | ✓ | 6K | Open Image V7 [19] | 9M |
| ScanObjectNN [20] | | 15K | Places [21] | 10M |
| 3D-Future[22] | ✓ | 16K | LSUN [23] | 59M |

Diffusion models, trained on billions of image-text pairs, have propelled the latest advancements in text-to-image generation, demonstrating the capability to produce high-fidelity images under textual prompts [24], [25], [26], [27], [28]. Utilizing pre-trained diffusion models for generating 3D digital content [29], [30] significantly reduces computational power requirements and dependence on 3D datasets, thereby greatly enhancing the feasibility and efficiency of 3D digital content generation. This paper meticulously investigates and analyzes

---

*Corresponding authors.

methods for generating 3D digital content, focusing on two key aspects: diffusion model priors and 3D representations. The generation of 3D digital content is categorized into two types based on the task: text-to-3D [29], [30], [31], [32], [33], [34], [35], [36], [37] and image-to-3D [35], [38], [39], [40], [41]. To compare the strengths and limitations of each approach, this study conducts a horizontal comparison of different models in terms of efficiency and quality. This paper also explores the challenges associated with generating 3D digital content using pre-trained diffusion models and discusses potential solutions to these issues.

Our contributions are summarized as follows:

- This paper delivers an exhaustive review and investigation of methods for generating 3D digital content, with a foundation in diffusion models.

- A horizontal comparison and analysis are conducted in this paper to discern variations in efficiency and quality among different models.

- Several viable solutions are proposed in this paper to address the current challenges in generating 3D digital content using diffusion models.

- Potential future research directions in the field of 3D digital content generation, guided by diffusion models, are outlined in this paper.

Additionally, it is worth noting that there is currently a lack of universally recognized evaluation metrics for text-to-3D digital content generation. We currently assess quality solely through visual observation, which introduces a certain level of subjectivity. In the realm of image-to-3D digital content generation, we will employ image-based metrics to objectively evaluate the generated 3D digital content. Furthermore, due to limitations in laboratory conditions, all experiments in this paper were conducted using a single A40 GPU, and the results are presented accordingly.

This paper is organized as follows: Section II introduces the relevant background knowledge on 3D representation methods and diffusion models. Section III conducts a comprehensive analysis and study of the schemes, algorithms, and workflows for both text-to-3D and image-to-3D conversions. Section IV provides a holistic evaluation of existing 3D content generation approaches, analyzing the strengths and limitations of different methodologies. Section V explores the current challenges and proposes envisioned solutions. Finally, the paper concludes with a summary and presents our thoughts on future research directions and themes in this field.

## II. RELATED WORK

The generation of 3D digital content based on diffusion models principally involves two components: 3D representation and diffusion priors. DreamFusion [29] pioneered the integration of diffusion models into the task of 3D digital content generation. Subsequent studies in this domain have been categorized into two approaches based on their characteristics: optimization-based methods [42] and multi-view prediction-based methods [43], [44]. The focal point of research in this field has been centered on optimizing 3D representations or fine-tuning diffusion models.

### A. 3D Representations

In the fields of computer graphics and computer vision, the 3D representation of objects encompasses various forms, including point clouds [45], [46], voxel grids [47], [48], meshes [49], [50], and implicit neural representations [51]. Each representation method has its distinct advantages and limitations, suitable for different types of 3D tasks. In research on 3D digital content generation based on diffusion models, Neural Radiance Fields (NeRF) [51] or 3D Gaussian Splatting [52] are commonly employed.

*1) Neural radiance fields:* NeRF uses a neural network to learn the continuous volume density and color of a scene [53]. Central to NeRF is the utilization of a Multi-Layer Perceptron to parametrically represent 3D objects, enabling high-quality synthesis of new viewpoint images. Theoretically, it can model shapes at any spatial resolution [54]. The MLP parameters, denoted as $\theta$, take the camera pose $c$ as input. The output comprises color and density. The process involves camera rays traversing the scene, generating a set of sample points along the ray path. The color and transparency of each sampled point on the ray are cumulatively processed to synthesize the color of each pixel. Subsequently, these colors and densities are utilized in volume rendering to generate the image $g(\theta, c)$. NeRF can learn from a series of 2D images taken from different angles and synthesize highly realistic new viewpoint images, which is crucial for achieving realistic 3D scene reconstruction.

*2) 3D gaussian splatting:* Structure-from-Motion (SfM) [55] can estimate point cloud distributions from a set of images using the COLMAP library. The work of 3D Gaussian Splatting starts with sparse SfM points, modeling the geometry as a set of 3D Gaussian functions. The fundamental idea of 3D Gaussian Splatting is to consider each point as the center of a Gaussian distribution. These points, rather than being isolated discrete entities, have a smooth, continuous weight distribution around them. Each point influences its surrounding area, quantified by a Gaussian function. Each 3D Gaussian is defined by the point's position, covariance matrix, and opacity $\alpha$. Specifically, the point's position is the mean of the 3D Gaussian, the covariance matrix determines the shape of the 3D Gaussian, and the opacity $\alpha$ is used for splatting, with spherical harmonics (SH) [56], [57] representing color. The method uses adaptive Gaussian densification to control the number and density of Gaussians per unit volume. This approach overcomes the issues of slow rendering speed or compromised image quality in previous methods, enabling high-quality, real-time novel view synthesis at 1080p resolution.

### B. Diffusion Models

Diffusion models consist of a forward process $q_{t, t \in [0,1]}$, and a reverse process $p_{t, t \in [0,1]}$. The forward process resembles a straightforward Brownian motion with time-varying coefficients [58]. Specifically, this process incrementally adds noise $\epsilon \in \mathcal{N}(0, \mathbf{I})$ to the original data $x_0$, thereby gradually transitioning the data distribution towards a Gaussian noise distribution [12], [59]. This step-by-step addition of noise effectively transforms the original data into a state that aligns with a predefined Gaussian distribution, laying the groundwork for the subsequent reverse process. Conversely, the reverse process employs a neural network to estimate the

noise added at each step of the forward process, progressively denoising the Gaussian distribution noise to ultimately restore the original data distribution. The distribution in the forward process is given by $q_t(x_t|x_0) := \mathcal{N}(\alpha_t x_0, \sigma_t^2 \mathbf{I})$ and $q_t(x_t) := \int q_t(x_t|x_0)q_0(x_0)\mathrm{d}x_0$. The coefficients $\alpha_t$ and $\sigma_t$ are selected to regulate the proportion of original data and noise. At the onset of the forward process, $\sigma_0 \approx 0$, while at the end, $\sigma_1 \approx 1$, where $\alpha_t^2 = 1 - \sigma_t^2$ [60], [61]. This careful adjustment of coefficients ensures a gradual and controlled transformation of the data. The reverse process, through a noise prediction network $\epsilon_\phi(x_t, t)$, predicts the noise added at each forward step. The overall training is conducted by minimizing

$$\mathcal{L}_{\text{Diff}}(\phi, x_0) = \mathbb{E}_{t,\epsilon}[\omega(t)||\epsilon_\phi(\alpha_t x_0 + \sigma_t \epsilon, t) - \epsilon||_2^2] \quad (1)$$

where, $\omega(t)$ is a weighting function that depends on the timestep $t$. the noise prediction network can be used for approximating the score function of both $q_t$ and $p_t$ by $S_\phi(x_t, t) = -\epsilon_\phi(x_t, t)/\sigma_t$.

Incorporating textual control within diffusion models enhances the controllability of the generated content [8]. Since each image adheres to a specific distribution pattern, utilizing the information embedded within the text as a directive allows for the progressive denoising of Gaussian noise images, culminating in the generation of images that align with the textual information. This process specifically involves training an encoder and a decoder, where the encoder maps images to a latent space and the decoder reconstructs images from this latent space data. The textual prompts $y$ are encoded using a text encoder $\tau_\theta(y)$ and are integrated into each step of the denoising process, which is trained by minimizing

$$\mathcal{L}_{\text{LDM}} = \mathbb{E}_{t,\epsilon,y}[\omega(t)||\epsilon_\phi(x_t, t, \tau_\theta(y)) - \epsilon||_2^2] \quad (2)$$

By introducing conditions into the noise reconstruction process, controlled image generation is achieved. This methodology exhibits robustness in producing high-resolution images with intricate details while maintaining the semantic structure of the images [62].

## III. METHODOLOGY

Diffusion models demonstrate extraordinary zero-shot capabilities in generating diverse images from textual descriptions. Fig. 1. demonstrates the ability of diffusion models to create multi-angular images using textual prompts.



Fig. 1. Generate multi-angle images based on text prompts.

Pre-trained diffusion models, having been trained with a vast array of internet data, have acquired an understanding of the distribution of images of most objects from various viewpoints [63]. By leveraging the geometric priors learned from natural images by large-scale diffusion models and integrating viewpoint control, fine-tuning these pre-trained models enables the generation of images from different perspectives. The viewpoint-conditioned diffusion models (Zero-1-to-3) [63] learn the relative control of camera perspectives using synthetic datasets, thereby facilitating the creation of novel views of the same object under specified camera transformations. Fig. 2. demonstrates the capability of the viewpoint-conditioned diffusion models to take a single-perspective image as input and generate images from diverse viewpoints.



Fig. 2. Generate different perspective images from a single viewpoint image.

The specific steps for using diffusion models as a prior to guide the generation of 3D digital content are as follows: First, initialize a 3D model, then continuously modify the shape of the 3D model according to the prompt. Upon completion of the iterative process, the final 3D model, when rendered from any perspective, aligns consistently with the content described in the prompt.

### A. Text-to-3D

The work on generating 3D digital content from textual prompts is built upon the foundations of text-to-image diffusion models [8], [26], [27], [28]. Given that the end product of diffusion models is an image, it's not feasible to directly use the results of diffusion models to supervise the generation of 3D digital content. However, it's possible to utilize the denoising process to guide this generation. The forward process of the diffusion model involves adding noise to the original data $x_0$ at timestep $t$, resulting in a noised image $\alpha_t x_0 + \sigma_t \epsilon$. During the reverse process, the noise prediction network estimates the noise $\epsilon$ added at each step, thus the denoised image can be represented as $x_\phi = [(\alpha_t x_0 + \sigma_t \epsilon) - \sigma_t \epsilon_\phi)]/\alpha_t$. This indicates that as long as the noise prediction is sufficiently accurate, the final image generated from Gaussian noise will also be accurate.

DreamFusion [29] employs NeRF as the 3D representation and utilizes a pre-trained text-to-image diffusion model as a critic. It achieves text-to-3D generation with impressive results through Score Distillation Sampling (SDS) loss. Specifically, the process involves rendering an image $x_{render}$ from a given viewpoint $c$ using the differentiable renderer $G(\theta, c)$. Here, $G$ is a differentiable rendering function parameterized by $\theta$, representing the parameters of the 3D object. Random amounts of noise are introduced into the rendered image $x_{render} := G(\theta, c)$ at various time steps $t$, resulting in $x_t = \alpha_t x_{render} +$

Fig. 3. A simplified framework for generating 3D digital content based on text prompts.

$\sigma_t \epsilon$ [30]. The pre-trained diffusion model predicts the sampling noise $\epsilon_\phi$ given a noisy image $x_t$, noise time step $t$, and text embedding $y$. It provides a gradient direction to update the 3D volumetric parameters $\theta$, with the overall gradient computed by the SDS function.

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\theta) = \mathbb{E}_{t,\epsilon,c}[\omega(t)(\epsilon_\phi(x_t,t,y) - \epsilon)\frac{\partial G(\theta,c)}{\partial \theta}] \quad (3)$$

Here, $\omega(t)$ is a weighting function. The scene model $G$ and the diffusion model $\phi$ can be considered as modular components. It can be demonstrated that this loss fundamentally measures the similarity between the rendered images and textual prompts [40]. During the iterative process, the SDS loss backpropagates only to update the NeRF parameters $\theta$, without altering the pre-trained diffusion model. As iterations progress, the 3D object gradually exhibits textures and geometric shapes that align with the textual prompt. The overall network architecture is succinctly illustrated in Fig. 3.

### B. Image-to-3D

People possess the ability to envision the 3D structure of an object from a single image, a skill largely derived from the vast amount of prior knowledge accumulated through life experiences. Much of the past research has focused on reconstructing 3D models from multi-angle images [56], [57], [64], [65]. This approach is intuitive, as multiple viewpoints are essential for acquiring 3D information. However, 3D reconstruction from multi-angle images remains inefficient. This method requires the collection and acquisition of images from multiple angles, implying that it can only reconstruct objects that already exist in the real world. An interesting aspect is that in industries with a high demand for 3D digital assets, such as gaming, virtual reality, and animation, the focus lies on innovative 3D models rather than mere reproductions of the real world. Typically, the creation of an original 3D model involves numerous steps, as illustrated in Fig. 4.



Fig. 4. 3D Model modeling process.

A promising approach to creating the requisite 3D models is through the generation of corresponding 3D models from a single image. While achieving controllability in diffusion models is a hot topic in further research [66], [67], [68], [69], there still lacks effective means to precisely control the images they generate. Consequently, the 3D models produced using text-to-3D methods may not always meet specific requirements. In other words, when inputting text prompts, no one can predict the structure of the 3D model until the result is generated. At this juncture, the task of image-to-3D conversion gains a significant advantage.

DreamFusion [29] achieves a text-to-3D generation method based on diffusion priors, demonstrating the exceptional capability of using diffusion priors to optimize NeRF. Related work [38], [41], [70] attempts to apply diffusion priors to single-image 3D generation. Owing to the fact that pre-trained diffusion models are primed with textual prompts, the approach for image-to-3D tasks diverges from that of text-to-3D tasks. Specifically, image-to-3D requires a process of textual inversion [68], differentiating it from the generation method used in text-to-3D tasks. A simplified network structure is illustrated in Fig. 5.

Fig. 5. A Two-stage framework for generating 3D digital content from a single image using diffusion priors.



Fig. 6. Qualitative comparisons of 3D digital content generation from textual descriptions.

The generation of 3D digital content from a single image is typically a two-stage process. The primary task of the coarse stage is to establish the model's basic outline, followed by refinement in the refine stage. Specifically, the coarse stage begins with preprocessing such as background removal [71], textual inversion [68], and depth estimation [72], [73] of the reference image. Background removal focuses on isolating the main object for modeling, while textual inversion generates

corresponding textual descriptions to guide the diffusion prior. Depth estimation provides a prior for depth information, supervising subsequent model generation. The overall process starts with initializing a 3D model, rendering images from random angles with added Gaussian noise, and then using a diffusion model to optimize the 3D model through back propagation using SDS loss and a series of reference image losses.

*1) Reference view reconstruction loss:* To ensure consistency between rendered images $G_\theta(c)$ from reference viewpoints $c$ and the reference images $x_0$ themselves, a reference view reconstruction loss is typically introduced at the reference viewpoints. This involves the use of Mean Squared Error (MSE) loss on the reference images and their masks.

$$\mathcal{L}_{rec} = \lambda_{rgb}||\mathrm{M} \odot (x_0 - G_\theta(c))||_2^2 + \lambda_{mask}||\mathrm{M} - M(G_\theta(c))||_2^2 \tag{4}$$

Here, $\theta$ represents the parameters of the 3D object being optimized, $\odot$ is Hadamard product, $\mathrm{M}$ is related to the mask, $M(\cdot)$ is the foreground mask acquired by the volume density along the ray of each pixel. $\lambda_{rgb}, \lambda_{mask}$ are the weights for the foreground RGB and the mask [38].

*2) Depth prior:* At reference viewpoints, relying solely on reference view reconstruction loss may result in poor geometric shapes. To address shape blur, indentations, and flatness, a depth prior is typically incorporated. Specifically, this involves using a pre-trained monocular depth estimator [72] to assess the depth $d$ of the reference image. The depth of the 3D content viewed from the reference viewpoint should closely match this depth prior. Generally, negative Pearson correlation is used for depth regularization.

$$\mathcal{L}_{depth} = -\frac{Cov(d(c), d)}{Var(d(c))Var(d)} \tag{5}$$

Here, $Cov(\cdot)$ denotes covariance, and $Var(\cdot)$ calculates standard deviation, $d(c)$ refers to the depth modeled at the reference viewpoint. Through the use of reference view reconstruction loss and depth prior loss, the alignment between the reference image and the 3D model at the reference viewpoint can be optimized as much as possible. Although the estimated depth may not accurately represent geometric details, it is sufficient to ensure reasonable geometric shape and resolve most ambiguities [40]. Furthermore, normal smoothness loss [38] and diffusion CLIP loss [40] can also be added.

*3) Diffusion prior:* The supervision of novel view generation is guided by a diffusion prior. Textual inversion is used to generate textual descriptions $y$ for the reference images. The SDS loss is employed for the continuous optimization of the 3D model.

The reference view loss includes details not captured by textual prompts, and SDS loss ensures the generated 3D model conforms to the object's expected shape. Combined, they ensure the model generation is faithful both to the reference image and to the textual prompts.

Upon completion of the coarse stage, the generated 3D model possesses a reasonable geometric shape, yet its overall geometric structure and texture remain somewhat rough. Based on the 3D model produced in the coarse stage, a model refinement network [76] can be utilized for further refinement,

enhancing its geometric structure and texture. The overall optimization process is fundamentally similar to that of the coarse stage.

## IV. EXPERIMENTS

In accordance with the primary research focus of this paper, we categorize the current frameworks for 3D digital content generation based on diffusion models into two distinct types: text-to-3D and image-to-3D. All experimental results were obtained using a single A40 GPU. Our analysis primarily concentrates on two key aspects: the quality of the generated content and the speed of generation.

### A. Text-to-3D

In the comparative experiments of text-to-3D digital content generation, we encountered frameworks that were either open-source or proprietary. For the open-source frameworks, experiments were conducted using the original codes from the respective papers. In the case of proprietary frameworks, we uniformly utilized threestudio [77] for experimentation. We acknowledge that there might be slight deviations in the results generated by threestudio compared to the original outcomes; however, we believe these differences do not significantly impact our evaluative conclusions. Additionally, in the realm of text-to-3D digital content generation, there are no universally accepted benchmarks for performance evaluation. Consequently, qualitative assessments were primarily based on visual inspections conducted by human observers. In our detailed experiments, we compare recent methods (DreamFusion [29], Latent-NeRF [74], Score Jacobian Chaining [34], ProlificDreamer [75], DreamGaussian [35]) for generating 3D objects from a textual prompt. Furthermore, considering the influence of textual prompt types on the model's generative performance, we employed two categories of textual descriptions: reality-based and imagination-based.The results of the generation are illustrated in Fig. 6.

Through a comparative analysis of the generated mesh quality and the overall generation time, as detailed in Table II, we observed that for objects existing in reality, ProlificDreamer [75] exhibits the highest quality of generation, albeit at the slowest speed. While DreamGaussian [35] may not match the former in terms of quality, it outperforms in generation speed. For imaginary objects, current mainstream frameworks struggle to achieve high-quality generation. We propose two avenues for optimization: firstly, refining textual prompts to more intricately describe the content envisioned, which could enhance the resultant generation. Secondly, augmenting the capabilities of the diffusion model by training it with larger datasets.

ProlificDreamer [75] proposed the use of Variational Score Distillation (VSD) to address issues such as over-saturation, over-smoothing, and low-diversity in the SDS loss. The core concept involves sampling within the distribution of 3D scenes, representing the 3D distribution with 3D parameter particles. A gradient-based particle updating rule is derived based on Wasserstein gradient flow. Despite its ability to achieve high-quality generation results, Prolificdreamer's method requires alternating training between LoRA [78] and NeRF during the training process, leading to prolonged training times. In contrast, DreamGaussian [35] employs 3D Gaussian Splatting [52]

TABLE II. MULTI-PERSPECTIVE COMPARATIVE ASSESSMENT OF TEXT-TO-3D DIGITAL CONTENT GENERATION FRAMEWORKS

| Method | DreamFusion [29] | Latent-NeRF [74] | Score Jacobian Chaining [34] | ProlificDreamer [75] | DreamGaussian [35] |
|---|---|---|---|---|---|
| 3D Representations | NeRF | NeRF | NeRF | NeRF | 3D Gaussian Splatting |
| Number of Stages | Single | Two | Single | Three | Two |
| Mesh Quality | ★ | ★★ | ★★★ | ★★★★★ | ★★★★ |
| Avg. Time | ∼40 minutes | ∼1 hour | ∼25 minutes | ∼13 hours | ∼4 minutes |

TABLE III. QUANTITATIVE RESULTS ARE PROVIDED FOR **PSNR ↑**, **LPIPS ↓**, AND **CLIP-SIMILARITY ↑**

| Dataset | Metrics | Zero-1-to-3 [63] | Magic123 [38] | DreamGaussian [35] | Stable Zero123 |
|---|---|---|---|---|---|
| RealFusion15 | PSNR↑ | 35.22 | 35.20 | **35.47** | 35.40 |
| | LPIPS↓ | 0.10 | 0.13 | 0.08 | **0.07** |
| | CLIP-Similarity↑ | 0.86 | **0.90** | 0.83 | 0.88 |

for 3D representation, significantly accelerating the generation speed.

### B. Image-to-3D

In the comparative experiments for image-to-3D digital content generation tasks, we utilized the RealFusion [41] dataset, comprising 15 distinct objects, for our analysis. We compare recent methods (Zero-1-to-3 [63], Magic123 [38], DreamGaussian [35], Stable Zero123) for generating 3D objects from a single unposed image, with specific experimental results depicted in Fig. 7. Unlike the generation of 3D digital content from textual prompt, the quality of 3D content generated from a single image can be assessed based on image-related metrics.

*1) PSNR:* PSNR is a widely used standard for quantifying the quality of image reconstruction or image compression. It measures the pixel-level differences between the original and the compressed or reconstructed image. PSNR is calculated based on the Mean Squared Error (MSE) between the two images. Generally, a higher PSNR value indicates that the reconstructed image is closer in quality to the original image. It primarily evaluates the pixel-level similarity between the reconstructed or compressed image and the original image, but it may not always align with human perceptual differences.

*2) LPIPS:* LPIPS is a more modern, deep learning-based metric used to assess the perceptual quality and similarity of images. LPIPS calculates the similarity by comparing the activations of a deep neural network when processing two images. This approach aims to more closely resemble the human visual perception system. LPIPS is used to evaluate the perceptual similarity of images, especially in cases where pixel-level metrics may not capture all aspects of human perception.

*3) CLIP-Similarity:* CLIP-Similarity is a metric used to evaluate the semantic similarity between images, based on features extracted by the CLIP model. Unlike traditional image similarity metrics that focus on pixel-level details, CLIP-Similarity measures how semantically or contextually similar two images are. CLIP-similarity is particularly useful when the evaluation criteria extend beyond mere visual or pixel-level accuracy and venture into the realm of contextual and conceptual alignment.

For evaluating the quality of generated 3D content from reference viewpoints, we follow the metrics used in previous studies [41], [70]. We employed PSNR and LPIPS [79] metrics to compare the rendered images against the reference images, thereby assessing the generation quality from reference viewpoints. For images rendered from novel viewpoints, the quality was evaluated using the CLIP-similarity [9], as presented in Table III. Moreover, because we preprocess the original images in the process of generating 3D digital content from images, we applied the same treatments to the rendered images of the final 3D models during comparisons to ensure the accuracy of experimental results.

Our findings reveal that DreamGaussian [35] exhibits the fastest generation speed and achieves the highest quality when viewed from a reference perspective. However, it is noteworthy that its performance in generating novel views is comparatively inferior. On the other hand, Magic123 [38] demonstrates superior performance in generating high-quality novel views by incorporating a dual prior in both 2D and 3D dimensions. Simultaneously, the experimental results also confirm that the combination of diffusion models and 3D Gaussian Splatting [52] can achieve rapid 3D digital content generation, although there is room for further improvement in generation quality.

## V. DISCUSSION

This study analyzes the frameworks related to text-to-3D content generation and image-to-3D content generation based on diffusion models, conducting extensive experiments. Through experimental comparative analysis, we identified numerous challenges in 3D content generation based on diffusion models.

### A. Current Issues

*1) Janus problem:* Due to the primary approach of utilizing the diffusion model to guide rendering images from various perspectives, subsequently directing the generation of 3D models, the Janus problem is pervasive in the task of 3D digital content generation based on the diffusion model.

*2) Over-Saturation:* Using SDS loss in the generation of 3D content leads to issues such as over-saturation, over-smoothing, and low-diversity problems.

*3) Controllability:* Achieving precise control in the generation of 3D content from text prompts is challenging, relying solely on textual cues.
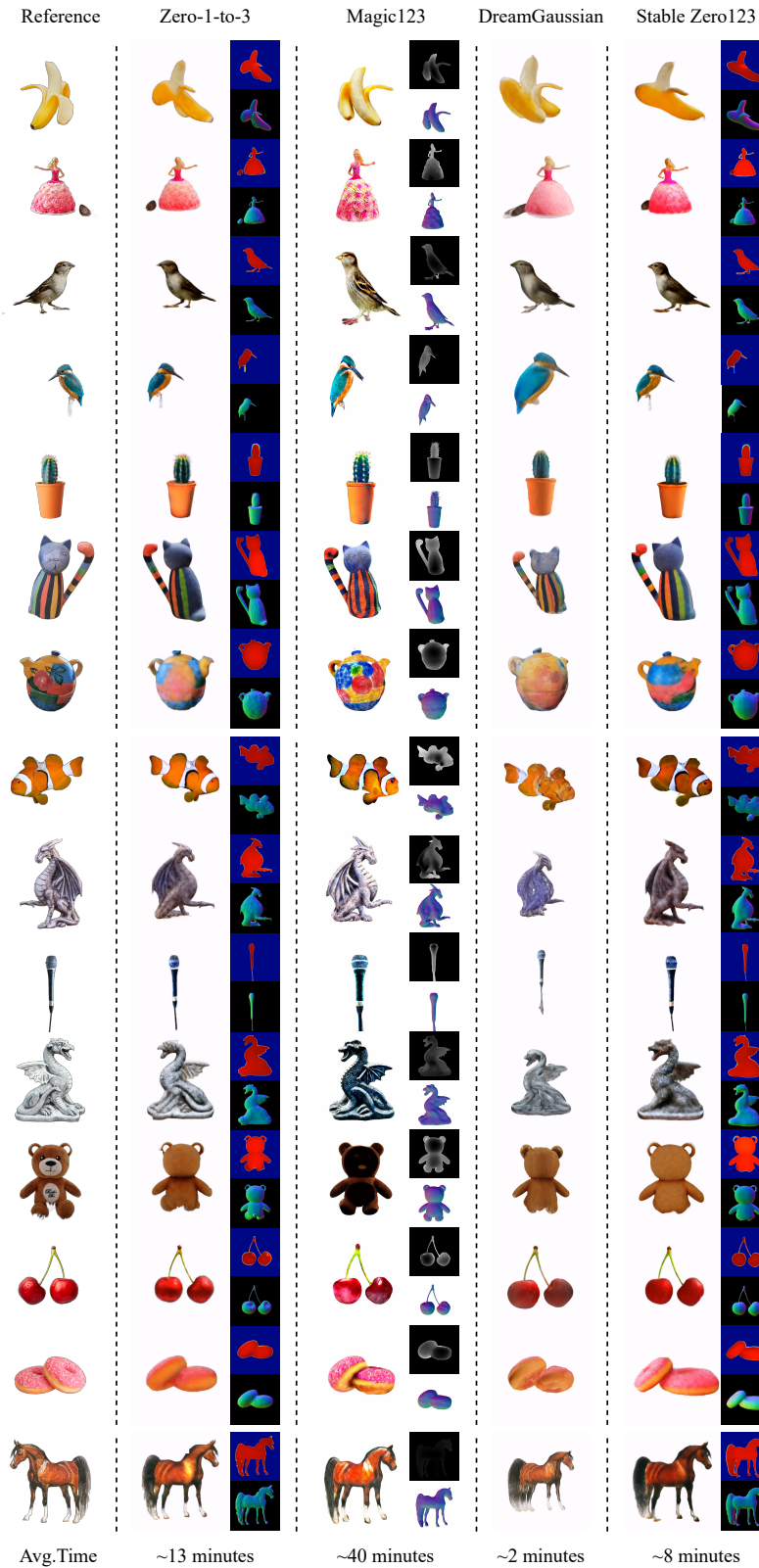
Fig. 7. Qualitative comparisons of 3D digital content generation from a single image.

*4) Editability:* Currently, there is no effective means to edit generated 3D content through artificial intelligence.

*5) Imagination:* Despite effective generation for real-world objects, the diffusion model struggles with the 3D reasoning

and imagination capabilities required for generating novel objects.

*6) Primary view dependency:* Tasks involving the generation of 3D content from a single image often require the input to be the primary view of the target object.

*7) Evaluation metrics:* A lack of a unified evaluation system for assessing the quality of generated 3D content.

*8) Generation quality:* Diffusion model-based 3D object generation faces issues of insufficient generation quality, resulting in objects that may lack realism or exhibit insufficient detail.

*9) Shape inconsistency:* Generated 3D objects may exhibit shape inconsistencies, particularly with complex geometric structures or topological relationships.

*10) Scale disparities:* Current 3D content generation models struggle to effectively handle objects of varying scales and are unable to generate 3D models of different sizes based on specific requirements.

### B. Potential Solutions to Some Issues

*1) Janus problem:* To address the Janus problem, employing multi-view [80] or 3D perception [37] diffusion models can help alleviate the issue. Additionally, an incremental modeling approach, similar to a "humanoid printer", can be applied, generating 3D models for partial views gradually.

*2) Over-Saturation:* An approach akin to that proposed by prolificdreamer [75], employing Variational Score Distillation (VSD), can be adopted to address the issue of over-saturation and further enhance the quality of generated 3D models. However, it is noteworthy that this method may lead to a reduction in efficiency.

*3) Controllability:* While achieving controllability in text-to-3D content generation tasks remains challenging, leveraging image-to-3D generation tasks can facilitate more controlled 3D content generation.

*4) Editability:* Editing of 3D content can be achieved through image editing techniques [66] or by combining Chat-GPT [1] to map text or voice into latent space for effective editing.

*5) Imagination:* In order to improve the generation performance of models, it is suggested to employ richer semantic description information. Alternatively, a more powerful diffusion model can be trained by incorporating a larger dataset. These strategies aim to enhance the overall effectiveness of the model in generating high-quality outputs.

*6) Primary view dependency:* Further enhancing the capabilities of novel view synthesis models to generate primary views of objects based on input images.

## VI. CONCLUSION

With the continuous development of generative artificial intelligence, the scope of generated content is expanding beyond text, audio, and image domains, gradually progressing towards the generation of 3D objects and environments. Fueled by the visions of virtual reality, augmented reality, and the metaverse, the demand for 3D digital content across various industries is expected to further burgeon.

Current research indicates that different frameworks for 3D digital content generation exhibit advantages and limitations in terms of both generation quality and efficiency. Through our specific investigations, we posit that the integration of diffusion models and 3D Gaussian Splatting will be a focal point in the future research of 3D digital content generation. Additionally, constrained by the controllability issue in text-to-3D, a viable workflow for 3D digital content generation is as follows: firstly, generate images from text, providing creators with creative input. Subsequently, employ artificial intelligence to optimize and edit the image content to achieve the desired appearance. Then, use an image-to-3D generation framework to create a 3D model. Finally, import the generated 3D model into 3D modeling software for further refinement.

With the advancement of 3D object generation frameworks, future research is expected to extend from individual objects to scene generation. How to integrate procedural scene generation with artificial intelligence in the future is a question worthy of consideration.

In summary, this review comprehensively elucidates how diffusion models can be leveraged for 3D digital content generation. We analyze key frameworks for 3D digital content generation and experimentally validate the efficiency and feasibility of combining diffusion models with 3D Gaussian Splatting for modeling. We summarize the existing challenges in 3D digital content generation based on diffusion models and propose potential solutions for some of these issues. Overall, we contend that image-to-3D digital content generation aligns more closely with societal applications, though we remain optimistic about the future of text-to-3D digital content generation.

## REFERENCES

[1] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

[2] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov, "Transformer-xl: Attentive language models beyond a fixed-length context," *arXiv preprint arXiv:1901.02860*, 2019.

[3] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.

[4] J. D. M.-W. C. Kenton and L. K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of naacL-HLT*, vol. 1, 2019, p. 2.

[5] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[6] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.

[7] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.

[8] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.

[9] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.

[10] J. Wu, C. Zhang, T. Xue, B. Freeman, and J. Tenenbaum, "Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling," *Advances in neural information processing systems*, vol. 29, 2016.

[11] J. Gao, T. Shen, Z. Wang, W. Chen, K. Yin, D. Li, O. Litany, Z. Gojcic, and S. Fidler, "Get3d: A generative model of high quality 3d textured shapes learned from images," in *Advances In Neural Information Processing Systems*, 2022.

[12] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *International conference on machine learning*. PMLR, 2015, pp. 2256–2265.

[13] Y. Song and S. Ermon, "Generative modeling by estimating gradients of the data distribution," *Advances in neural information processing systems*, vol. 32, 2019.

[14] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.

[15] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[16] L. Liu, W. Xu, H. Fu, S. Qian, Q. Yu, Y. Han, and C. Lu, "Akb-48: A real-world articulated object knowledge base," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 14 809–14 818.

[17] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.

[18] T. Wu, J. Zhang, X. Fu, Y. Wang, J. Ren, L. Pan, W. Wu, L. Yang, J. Wang, C. Qian *et al.*, "Omniobject3d: Large-vocabulary 3d object dataset for realistic perception, reconstruction and generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 803–814.

[19] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Malloci, A. Kolesnikov *et al.*, "The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale," *International Journal of Computer Vision*, vol. 128, no. 7, pp. 1956–1981, 2020.

[20] M. A. Uy, Q.-H. Pham, B.-S. Hua, T. Nguyen, and S.-K. Yeung, "Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1588–1597.

[21] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 6, pp. 1452–1464, 2017.

[22] H. Fu, R. Jia, L. Gao, M. Gong, B. Zhao, S. Maybank, and D. Tao, "3d-future: 3d furniture shape with texture," *International Journal of Computer Vision*, vol. 129, pp. 3313–3337, 2021.

[23] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao, "Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop," *arXiv preprint arXiv:1506.03365*, 2015.

[24] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," *arXiv preprint arXiv:2112.10741*, 2021.

[25] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8821–8831.

[26] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36 479–36 494, 2022.

[27] Y. Balaji, S. Nah, X. Huang, A. Vahdat, J. Song, K. Kreis, M. Aittala, T. Aila, S. Laine, B. Catanzaro *et al.*, "ediffi: Text-to-image diffusion models with an ensemble of expert denoisers," *arXiv preprint arXiv:2211.01324*, 2022.

[28] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint arXiv:2204.06125*, vol. 1, no. 2, p. 3, 2022.

[29] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, 2022.

[30] C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu, and T.-Y. Lin, "Magic3d: High-resolution text-to-3d content creation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 300–309.

[31] A. Jain, B. Mildenhall, J. T. Barron, P. Abbeel, and B. Poole, "Zero-shot text-guided object generation with dream fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 867–876.

[32] N. Mohammad Khalid, T. Xie, E. Belilovsky, and T. Popa, "Clip-mesh: Generating textured meshes from text using pretrained image-text models," in *SIGGRAPH Asia 2022 conference papers*, 2022, pp. 1–8.

[33] R. Chen, Y. Chen, N. Jiao, and K. Jia, "Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation," *arXiv preprint arXiv:2303.13873*, 2023.

[34] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich, "Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 619–12 629.

[35] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng, "Dreamgaussian: Generative gaussian splatting for efficient 3d content creation," *arXiv preprint arXiv:2309.16653*, 2023.

[36] W. Li, R. Chen, X. Chen, and P. Tan, "Sweetdreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d," *arXiv preprint arXiv:2310.02596*, 2023.

[37] J. Sun, B. Zhang, R. Shao, L. Wang, W. Liu, Z. Xie, and Y. Liu, "Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior," *arXiv preprint arXiv:2310.16818*, 2023.

[38] G. Qian, J. Mai, A. Hamdi, J. Ren, A. Siarohin, B. Li, H.-Y. Lee, I. Skorokhodov, P. Wonka, S. Tulyakov *et al.*, "Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors," *arXiv preprint arXiv:2306.17843*, 2023.

[39] X. Long, Y.-C. Guo, C. Lin, Y. Liu, Z. Dou, L. Liu, Y. Ma, S.-H. Zhang, M. Habermann, C. Theobalt *et al.*, "Wonder3d: Single image to 3d using cross-domain diffusion," *arXiv preprint arXiv:2310.15008*, 2023.

[40] J. Tang, T. Wang, B. Zhang, T. Zhang, R. Yi, L. Ma, and D. Chen, "Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior," *arXiv preprint arXiv:2303.14184*, 2023.

[41] L. Melas-Kyriazi, I. Laina, C. Rupprecht, and A. Vedaldi, "Realfusion: 360deg reconstruction of any object from a single image," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 8446–8455.

[42] Z. Chen, F. Wang, and H. Liu, "Text-to-3d using gaussian splatting," *arXiv preprint arXiv:2309.16585*, 2023.

[43] Y. Liu, C. Lin, Z. Zeng, X. Long, L. Liu, T. Komura, and W. Wang, "Syncdreamer: Generating multiview-consistent images from a single-view image," *arXiv preprint arXiv:2309.03453*, 2023.

[44] H. Weng, T. Yang, J. Wang, Y. Li, T. Zhang, C. Chen, and L. Zhang, "Consistent123: Improve consistency for one image to 3d object synthesis," *arXiv preprint arXiv:2310.08092*, 2023.

[45] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.

[46] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in neural information processing systems*, vol. 30, 2017.

[47] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1912–1920.

[48] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2015, pp. 922–928.

[49] A. Sinha, A. Unmesh, Q. Huang, and K. Ramani, "Surfnet: Generating 3d shape surfaces using deep residual networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 6040–6049.

[50] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang, "Pixel2mesh: Generating 3d mesh models from single rgb images," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 52–67.

[51] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[52] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics (ToG)*, vol. 42, no. 4, pp. 1–14, 2023.

[53] K. Gao, Y. Gao, H. He, D. Lu, L. Xu, and J. Li, "Nerf: Neural radiance field in 3d vision, a comprehensive review," *arXiv preprint arXiv:2210.00379*, 2022.

[54] Z. Shi, S. Peng, Y. Xu, A. Geiger, Y. Liao, and Y. Shen, "Deep generative models on 3d representations: A survey," *arXiv preprint arXiv:2210.15663*, 2022.

[55] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.

[56] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa, "Plenoxels: Radiance fields without neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 5501–5510.

[57] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.

[58] T. Karras, M. Aittala, T. Aila, and S. Laine, "Elucidating the design space of diffusion-based generative models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 26 565–26 577, 2022.

[59] H. Cao, C. Tan, Z. Gao, Y. Xu, G. Chen, P.-A. Heng, and S. Z. Li, "A survey on generative diffusion model," *arXiv preprint arXiv:2209.02646*, 2022.

[60] D. Kingma, T. Salimans, B. Poole, and J. Ho, "Variational diffusion models," *Advances in neural information processing systems*, vol. 34, pp. 21 696–21 707, 2021.

[61] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020.

[62] J. Li, Z. Li, Y. Li, and L. Wang, "P-2.12: A comprehensive study of content generation using diffusion model," in *SID Symposium Digest of Technical Papers*, vol. 54. Wiley Online Library, 2023, pp. 522–524.

[63] R. Liu, R. Wu, B. Van Hoorick, P. Tokmakov, S. Zakharov, and C. Vondrick, "Zero-1-to-3: Zero-shot one image to 3d object," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9298–9309.

[64] Y. Wang, Q. Han, M. Habermann, K. Daniilidis, C. Theobalt, and L. Liu, "Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3295–3306.

[65] C. Reiser, R. Szeliski, D. Verbin, P. Srinivasan, B. Mildenhall, A. Geiger, J. Barron, and P. Hedman, "Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes," *ACM Transactions on Graphics (TOG)*, vol. 42, no. 4, pp. 1–12, 2023.

[66] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 3836–3847.

[67] O. Avrahami, D. Lischinski, and O. Fried, "Blended diffusion for text-driven editing of natural images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18 208–18 218.

[68] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or, "An image is worth one word: Personalizing text-to-image generation using textual inversion," *arXiv preprint arXiv:2208.01618*, 2022.

[69] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 500–22 510.

[70] D. Xu, Y. Jiang, P. Wang, Z. Fan, Y. Wang, and Z. Wang, "Neurallift-360: Lifting an in-the-wild 2d photo to a 3d object with 360deg views," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4479–4489.

[71] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, "Segment anything," *arXiv:2304.02643*, 2023.

[72] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun, "Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer," *IEEE transactions on pattern analysis and machine intelligence*, vol. 44, no. 3, pp. 1623–1637, 2020.

[73] S. M. H. Miangoleh, S. Dille, L. Mai, S. Paris, and Y. Aksoy, "Boosting monocular depth estimation models to high-resolution via content-adaptive multi-resolution merging," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9685–9694.

[74] G. Metzer, E. Richardson, O. Patashnik, R. Giryes, and D. Cohen-Or, "Latent-nerf for shape-guided generation of 3d shapes and textures," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 663–12 673.

[75] Z. Wang, C. Lu, Y. Wang, F. Bao, C. Li, H. Su, and J. Zhu, "Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation," *arXiv preprint arXiv:2305.16213*, 2023.

[76] T. Shen, J. Gao, K. Yin, M.-Y. Liu, and S. Fidler, "Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis," *Advances in Neural Information Processing Systems*, vol. 34, pp. 6087–6101, 2021.

[77] Y.-C. Guo, Y.-T. Liu, R. Shao, C. Laforte, V. Voleti, G. Luo, C.-H. Chen, Z.-X. Zou, C. Wang, Y.-P. Cao, and S.-H. Zhang, "threestudio: A unified framework for 3d content generation," https://github.com/threestudio-project/threestudio, 2023.

[78] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021.

[79] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[80] Y. Shi, P. Wang, J. Ye, M. Long, K. Li, and X. Yang, "Mvdream: Multi-view diffusion for 3d generation," *arXiv preprint arXiv:2308.16512*, 2023.

# A Robust Deep Learning Model for Terrain Slope Estimation

Abdulaziz Alorf

Department of Electrical Engineering, College of Engineering,

Qassim University, Buraydah 52571, Saudi Arabia

*Abstract*—Interest in autonomous robots has grown significantly in recent years, motivated by the many advances in computational power and artificial intelligence. Space probes landing on extra-terrestrial celestial bodies, as well as vertical take-off and landing on unknown terrains, are two examples of high levels of autonomy being pursued. These robots must be endowed with the capability to evaluate the suitability of a given portion of terrain to perform the final touchdown. In these scenarios, the slope of the terrain where a lander is about to touch the ground is crucial for a safe landing. The capability to measure the slope of the terrain underneath the vehicle is essential to perform missions where landing on unknown terrain is desired. This work attempts to develop algorithms to assess the slope of the terrain below a vehicle using monocular images in the visible spectrum. A lander takes these images with a camera pointing in the landing direction at the final descent before the touchdown. The algorithms are based on convolutional neural networks, which classify the perceived slope into discrete bins. To this end, three convolutional neural networks were trained using images taken from multiple types of surfaces, extracting features that indicate the existing inclination in the photographed surface. The metrics of the experiments show that it is feasible to identify the inclination of surfaces, along with their respective orientations. Our overall aim is that if a hazardous slope is detected, the vehicle can abort the landing and search for another, more appropriate site.

*Keywords*—*Terrain slope estimation; spacecrafts; robotics; artificial intelligence; machine learning techniques; deep neural network; computer vision*

## I. INTRODUCTION

In recent decades, interest in autonomous robots has significantly risen. Through the use of autonomous systems, we refer to systems that are conceived with capabilities to make certain types of decisions during the execution of their missions, minimizing the need for human-operator interventions. Their capability to accomplish missions that might be hazardous for humans (e.g., exploring celestial bodies or disaster monitoring), as well as to automate everyday tasks (e.g., driving cars or moving objects within a warehouse), make them highly valuable. They are currently the object of multiple research efforts addressing a wide range of challenges.

Space probes navigating in the proximity and landing on extraterrestrial celestial bodies [1], [2], [3] or indoor [4] and outdoor [5] navigating drones are among the most common applications of autonomous systems. For a robot intended to land on an asteroid (or other unexplored celestial bodies), the details of the terrain and local slope of its surface are very likely unknown, unless the celestial body had been previously studied thoroughly (e.g., the Moon or Mars) [6]. Moreover,

as their missions might take place at large distances from Earth, real-time communication with ground stations might not be feasible, as the signals would require on the order of a number of minutes to traverse the path from the probe to Earth. Therefore, capabilities for autonomous navigation and landing are highly required. Likewise, vertical take-off and landing (VTOL) aircraft exploring unknown terrains also constitute highly demanded applications of these systems [7], [8]. In the aforementioned scenarios, landing on surfaces that are not well-known or mapped would require the lander to make decisions on how to approach and where to touch ground, thus minimizing the intervention of human operators [6], [9]. Fig. 1 shows an area of the Martian surface, with its many geographical features, including flat and inclined areas. A completely autonomous robot approaching the surface for landing should be able to determine whether it is going to touch down on a flat or inclined area.



Fig. 1. Photograph of a piece of the martian surface. Credit: NASA/JPL-Caltech/Univ. of Arizona.

When a space probe or VTOL aircraft lands, the final descent to the surface intends to follow position and attitude trajectories that ensure a smooth touchdown, with the entire landing gear leaning simultaneously on the ground. In general, these trajectories aim to touch ground at areas where a vector normal to the local terrain is somewhat aligned with the negative local gravity vector [6], [10]; however, there have been some efforts to study the feasibility of take-off and landing from sloped terrain [11]. Fig. 2 illustrates this concept.

In Fig. 2, there are two landers, one denoted by A (on the left hand-side of the figure), and the second denoted by B (at the center of the figure). The arrows labeled as DM show the

Fig. 2. Schematic view of robots landing on flat and inclined terrain.

direction of motion of both landers, and the local gravity vector is also depicted. The two dextral coordinate systems, denoted as $B$, fixed to the bodies of the landers A and B, respectively, can also be observed. The unit vectors $z_B$ always point along the legs of the body, while the unit vectors $y_B$ point towards the right side of it, as can be seen in Fig. 2. Correspondingly, $x_B = y_B \times z_B$.

A hypothetical piece of terrain where the probes would land is also displayed in Fig. 2. Assume that the piece of terrain consists of multiple piecewise planar surfaces, $k$, that are continuously concatenated, and each of which has a different slope represented by a local normal unit vector $N_k$ associated with it. The depicted scene shows Lander A about to touch ground on an area with $z_B = -N_1$. In contrast, Lander B is about to touch ground on a highly inclined region with respect to $z_B$. This may cause the lander to touch ground on its left leg. If the lander ignores the inclination of the surface upon touching ground, it might constitute a hazard for the landing moment. The robot would touch ground with only its left leg, perhaps overloading it; alternatively, it could turn over, roll down, or even slip down, leading to an undesired termination of the mission. For sake of clarity, $z_{B\_A}$ and $z_{B\_B}$ will hereafter be indistinctly referred to as $z_B$.

In order to avoid these hazards, these robots should be endowed with the capability to sense that the surface underneath is highly inclined, with respect to its attitude (as with Surface 2); and, hence, that it is necessary to navigate to a more suitable location for touching ground (as with Surface 1). Understanding information on the terrain and the slope of

the surface underneath constitutes an essential task in any of the following cases: (I) A team at the control center makes the decision on where and how to land, and imparts the commands to the robot; or (II) the robot decides by itself where to land and how to accomplish it.

Current and past missions have used distance-based sensors, such as a radar or lidar [12], [13], [14], in order to evaluate the slopes of the terrain underneath. In both cases, electromagnetic waves are sent towards an object or surface, with respect to which the measuring distance is determined. The time at which the reflection of these waves reaches the sensor is measured, which enables the computation of the distances to multiple points. An advantage of these sensors is the accuracy with which they can measure the terrain; however, they are usually expensive and power-consuming. Motivated by the aforementioned scenarios, we explore the feasibility of using individual images from a monocular camera (in the visible spectrum) to classify the relative inclination between the focal plane of the camera (associated with the vector $z_B$ in Fig. 2) and the piece of terrain photographed by the camera (associated with the vector $N$ in Fig. 2). This contribution would allow a lander to use monocular images to evaluate the inclination of the terrain underneath and assess whether or not it is a good site to touch ground. This work proposes an alternative methodology to evaluate the slope of the terrain, based on the processing of monocular images with artificial neural networks. The aim of this concept is not to substitute the usage of radars and lidars, but to assess another operational principle that could be used on its own

or fused with information retrieved from the aforementioned sensors. To the best of our knowledge, we could not find approaches similar to the one presented herein, which is why we considered it interesting to evaluate the feasibility of the approaches elaborated in the following.

### A. Related Work

Computer vision applications for space vehicle navigation have been pursued intensively in recent decades [15]. In this respect, photo-cameras also constitute sensors that can be used for relative attitude determination and position. When two or more spacecraft are flying in proximity or performing docking maneuvers, images in the visible spectrum can provide highly relevant information to determine the relative position and attitude between them [9], [16], [17]. Images or videos recorded by the cameras are processed by algorithms that identify feature points (which might be pre-defined or not) and track them along the sequence of images. The position of these feature points in the images provides an indication about the relative position and orientation between the two spacecraft during the maneuvers. In these types of applications, the information retrieved from images is usually fused with information provided by other sensors, such as gyros, range finders, and star trackers.

In the context of aerial vehicles, the use of computer vision techniques has also been intensively explored, especially for landing and collision avoidance purposes. In [18] and [19], the usage of optical flow measurements for deriving control laws for landing VTOL unmanned aerial vehicles on moving platforms was investigated. In study [20], optical flow measurements were exploited for the determination of landing control laws of UAVs in cluttered environments. In study [21] and [22], optical flow was also used, but for collision avoidance purposes. In study [23], neural networks were implemented to process video signals for indoor navigation assistance in maneuver planning. In study [24], a thorough survey of the vision-based techniques used for UAV navigation was presented.

There have also been efforts to determine terrain slopes from multiple overlapping aerial or satellite imagery [25], [26], [27], [28], [29], [30], [31]; however, these articles were not specifically intended for implementations in autonomous systems that make decisions in real-time from individual images, as they need to combine multiple overlapping images to determine the slope of the photographed terrain.

Extracting 3D features from 2D scene projections (images) has been considered a central challenge in computer vision. It has been extensively tackled in a diversity of contexts, and through multiple approaches. Among them, using texture cues to understand the shapes of 3D geometries has been highly pursued [32], [33], [34]. In research [35], the recovery of 3D shapes from the observed distortions in the density of the textures was analyzed, and the corresponding equation for determining the shapes of planar and curved surfaces was derived. In study [36], [37], and [38], affine transforms were proposed to model the relationship between the texture distortion and direction of the points in the image. Using these transforms, the orientation and shape parameters were estimated for multiple directions in the image. In study [39], the authors adopted 3D morphable models that are fitted to

pixel intensity, edges, and specular highlights, thus maximizing the posterior probability of the parameters upon the input image.

The estimation of depth (the coordinate along the line of sight) from monocular images also constitutes another highly pursued technical challenge, as local features do not represent enough information to estimate the depth of arbitrary points in images. Usually, depth is perceived as a result of two or more associated vision sensors (stereo-vision). In study [40], a Markov Random Field was used to learn depth cues from monocular images, which were used to reinforce a stereo vision system. In study [41], depth maps for still scenes were estimated, based on depth cues captured by supervised learning algorithms. In study [42], using the Lucas–Kanade method, the authors aimed to estimate the depths in scenes captured by a moving camera. In research [43], the authors exploited de-focus (blur) and textures to estimate depth maps. In study [44], a CNN was utilized to estimate the relative and absolute depth maps, which were optimally combined. In study [45], a ranking approach was used for relative depth estimation.

Terrain characterization can be considered a specific application of 3D shape estimation. In robotics navigation, terrain characterization is essential for the landing and ground traversal of robots performing on unknown terrains. Computer vision applications have also been extensively used for this purpose. In study [46], the authors proposed a system for terrain recovery which could be used for autonomous rotorcraft. The system aims to recover the material properties and geometry of the local terrain, using stereo cameras, a global navigation satellite system (GNSS), and inertial measurement units (IMU). In study [47], principal component analysis was used to estimate the normal vectors of surfaces, in order to determine traversable and non-traversable areas from depth images. This method was implemented in a rover-like robot, equipped with cameras and depth sensors, which allowed for the measurement of three-dimensional point clouds. In researches [48] and [49], algorithms based on convolutional neural networks (CNN) were presented for the classification of different terrain types and surface features, from both orbital and ground images. These outcomes were then used to determine terrain traversability for rovers. Terrain classification has been also pursued based on proprioceptive signals [50].

This work attempts to use convolutional networks to extract features from monocular images, thus allowing for estimation of the slope of the piece of terrain just below a landing vehicle. We presumed that a properly trained neural network can find, in monocular images, indicators of the slopes and orientation of the terrain below. These cues would mainly be in the variations of the visual textures (densities and sizes) across an image. To further define the problem tackled in this article, several more concepts are introduced.

### B. Incidence Angle, Far and Near Sides

In many scenarios where the human eye is looking at a given surface, the brain can understand whether the surface is normal or inclined with respect to the line-of-sight (LOS)[1].

---

[1]By line-of-sight, we refer to the line joining our eyes (or the camera) with the aimed object

This is an ability that we learn when our vision system commences its development, which improves as we are exposed to more types of surfaces with different inclinations. We do not have the capability to accurately measure the relative angle between the LOS and the plane of the surface we are looking at. However, we can distinguish qualitatively high relative inclinations from low or null inclinations. For instance, consider Fig. 3a and Fig. 3b, which show images of a street surface paved with stones, taken at different angles. $i$ denotes the incident angle, which is defined as the angle between the LOS and the vector $N$ normal to a given surface.



(a) Incident angle $i \simeq 0$.



(b) Incident angle $> 0$.

Fig. 3. Surfaces perceived as looking with different incident angles.

Fig. 4a and Fig. 4b show a person looking at a surface from $i = 0$ and $i > 0$ deg, respectively. Generally speaking, a person could determine, only from an image, whether $i > 0$ or not. The reader could probably determine that the surface shown in Fig. 3b has a higher incident angle, with respect to the surface of the street, than that of Fig. 3a. Clearly, the texture of the surface we are looking at helps us to distinguish the angle of incidence. A purely plain surface with an incident angle would be impossible to distinguish from the same surface with a null incident angle.

Furthermore, it can also be seen, from Fig. 3b, that the upper part of the image was at a larger (far side) distance



(a) Null incident angle, $i = 0$.



(b) Incident angle $> 0$.

Fig. 4. Person looking at a surface from different incident angles.

from the focal plane of the camera than the lower part (near side) of it. The *far side* and *near side* of an image denote the sides of it that are farther from and closer to the focal plane, respectively. This is illustrated in Fig. 5a and Fig. 5b. Fig. 5a illustrates the case where the angle of incidence of the camera with respect to the surface is null, which would correspond to the Lander A of Fig. 2 if it had a camera pointing along its $z_B$ axis. Fig. 5b represents a scenario where the angle of incidence is considerably greater than zero, which would correspond to what would be seen by a camera pointing along $z_B$ in Lander B. In the latter case, the far and near sides are indicated in the image.

Moreover, for images characterized by $i > 0$, we define a roll angle, $r$, as the angle measured clockwise between a vector from the center of the image pointing towards the far side of it and a vector from the center of the image pointing towards the upper edge of it.

With the definitions stated above, the goal of this article is to report two experiments, aimed at deriving algorithms that

(a) Null incident angle, $i = 0$.



(b) Incident angle, $i > 0$.

Fig. 5. Difference in geometries for images taken with $i = 0$ (case A) and with $i > 0$ (case B).

solve the following problem: *With a single monocular image of a surface in the visible spectrum, classify the angle between the normal vector at the photographed surface and the optical axis of the camera (i.e., $z_B$) using the following categories:*

- *Normal* (no inclination), $r$ undefined, denoted as type N, as seen in Fig. 3a;

- *Upward* inclination, $r = 0$ deg, denoted as type UI, as seen in Fig. 3b;

- *Downward* inclination, $r = 180$ deg, denoted as type DI, as seen in Fig. 6a;

- *Leftward* inclination, $r = 270$ deg, denoted as type LI, as seen in Fig. 6b;

- *Rightward* inclination, $r = 90$ deg, denoted as type RI, as seen in Fig. 6c.



(a) Surface with incident angle $i > 0$ and $r = 180$ deg.



(b) Surface with incident angle $i > 0$ and $r = 270$ deg.



(c) Surface with incident angle $i > 0$ and $r = 90$ deg.

Fig. 6. Different angles $r$ for images with $i > 0$ deg.

A few observations follow. Tackling the problem as a

classification problem constitutes a coarser measurement and is fairly easier than determining the angles precisely as continuous real numbers. However, it is an essential first step in the major goal of an ongoing project, aimed at precisely estimating these angles as real continuous variables. Indeed, one of the described experiments classifies the images considering three ranges of values for the incidence angle: 0 deg, 20 deg, and 40 deg. The two aforementioned experiments are elaborated in the following sections.

In order to derive a solution to the problem stated above, we trained convolutional neural networks using images of diverse types of surfaces. These surfaces represent ground terrains where a lander might perform a touchdown. The appearance of a landing surface can be very diverse. Therefore, we intended to obtain an algorithm that can work well with multiple textures, including some that can be found in the Earth's ground types and others that are not necessarily observed as ground surfaces. Including images with multiple textures allows the CNNs to learn features that can provide information about the angles $i$ and $r$, even when the angle is sensed across images having different textures. Section II-C displays samples of the types of surfaces used in this work.

Artificial neural networks (ANNs) have been used extensively for classification problems [51]. They represent algorithms that can be very efficient for solving high-dimensional classification problems. Once they have been trained, they can process given inputs and return the probabilities that these inputs to belong to any of the classes for which they have been trained. Training refers to the process that optimizes the internal coefficients of the operator to the specific classification process intended [52], [53].

In the field of image processing, CNNs constitute a powerful tool for image classification and interpretation [54]. CNNs are special types of NN that are well-suited to dealing with images. The images are introduced as tensors, where each entry of the tensor dictates the color intensity of its corresponding pixel. Upon every input image, the CNN performs a sequence of mathematical operations on the numerical values assigned to each pixel, and computes the probability of the whole image (or a part/parts of it) corresponding to a given pre-defined class.

CNNs are mainly composed of convolutional layers and, possibly, other types of intermediate layers. A convolutional layer is a portion of the algorithm where convolutional filters are applied to the input of that layer. By means of these convolutional filters, CNNs perform convolutional operations on these images, extracting features that provide indications about which of the pre-defined classes best characterize the image analyzed [55].

Typical applications of CNNs include object recognition within images [56], determining the position of sought objects within images [57], [58], and identifying pathologies such as CoViD-19 in pulmonary radiographies [59], [60], [61]. CNNs can have multiple architectures. The architecture of a CNN refers to the sequence in which the multiple operations are applied to the input. In a given classification problem, different architectures can produce different results, with the same inputs. In this work, we intend to exploit these classification tools to determine which of the aforementioned classes would best characterize an image, thereby associating each image

with the most suitable ranges for $i$ and $r$.

The approach proposed herein should enable a lander robot to interpret how appropriate for landing, in terms of inclination, the piece of terrain underneath it is. To date, CNNs have not been used to analyze angles between the focal plane and the photographed surface. The contribution of this work consists of a methodology to create algorithms that can classify the incidence angle of an image into certain pre-defined categories. These algorithms could become an essential tool for space probes or autonomous VTOL aircraft, in order to determine whether a piece of terrain has a slope that makes it appropriate landing site.

The remainder of the paper is structured as follows: Section II describes the image collection and preparation processes, as well as the architecture of the convolutional neural network implemented for the classifier. Section III elaborates on the obtained results and discusses potential directions for improvement. Finally, Section IV presents our main highlights and observations in this work.

## II. METHODOLOGY

This work reports two experiments in which, using CNNs, we address two versions of the problem stated in Section I. These experiments are referred to as Experiment I and Experiment II, and are described in the following. The proposed pipeline, which represents the procedures followed by the experiments, is shown in Fig. 7. The figure shows the process of collecting images, augmenting the data set, and training and testing the respective convolutional neural networks.

### A. Experiment I

This experiment aimed to derive algorithms to solve the following problem: *With a single image of a surface, taken by a single camera in the visible spectrum, classify the angle between the normal vector at the photographed surface and the optical axis of the camera (i.e., $z_B$) into the following categories*:

- *Normal* (no inclination), $r$ undefined, denoted as type N, as seen in Fig. 3a;

- *Upward* inclination, $r = 0$ deg, denoted as type UI, as seen in Fig. 3b;

- *Downward* inclination, $r = 180$ deg, denoted as type DI, as seen in Fig. 6a;

- *Leftward* inclination, $r = 270$ deg, denoted as type LI, as seen in Fig. 6b;

- *Rightward* inclination, $r = 90$ deg, denoted as type RI, as seen in Fig. 6c.

As previously stated, these problems were tackled using CNNs. At present, there exist several CNN architectures that are renowned for having very good performance in certain types of image classification problems. To address this experiment, we implemented two different architectures and compared the exhibited performance. These architectures are described in Section II-A1 and Section II-A2.

Fig. 7. Schematic representation of the pipeline which represents the procedures followed by the experiments.

*1) Convolutional neural network based on the VGG16 architecture:* First, we implemented a convolutional neural network based on the VGG16 architecture. This architecture was proposed by the Visual Geometry Group at the University of Oxford, winning the ILSVR (Imagenet) competition in 2014 [56]. Since then, it has become widely used, thanks to its excellent performance. It was described in the seminal work by Simonyan and Zisserman [62] and, since then, has been implemented in a variety of applications [63], [64], [65]. A schematic view of the architecture of this CNN is presented in Fig. 8.

In its original implementation, the input of the first convolutional layer had a fixed size of $224 \times 224$ and was of RGB type. The image was processed through an array of 13 successive convolutional layers with max pooling layers in between every two or three convolutional layers, as displayed in Fig. 8. Convolutional filters of $3 \times 3$ were used in every convolutional layer, with stride of 1 pixel and padding of 1 pixel. Spatial pooling was performed in a total of five

max-pooling layers. The max-pooling was performed through windows of $2 \times 2$ pixels and a stride of 2 pixels. After the convolutional layers, three dense layers with 4096 nodes each were used, followed by the final softmax layer. The hidden layers were built using rectified linear unit (ReLU) activation functions.

In this work, the implementation of the VGG16 network was accomplished using the Keras API [66]. Keras is a programming interface for Tensorflow, which is a Google open-source library for machine learning applications [67]. Keras has a built-in implementation of the VGG16 network that allows for image sizes different from the original configuration of $224 \times 224$ pixels, enabling us to define a model using an image size of $300 \times 300$. Keras also enables the user to remove the three fully connected layers at the end of the network that the original version of VGG16 had. In this work, these layers were removed, and the output of the convolutional layers was passed through a filter of the average pool type with a size of $2 \times 2$ pixels. Subsequently, a single fully connected layer of

Fig. 8. Schematic diagram of the VGG16 CNN architecture. Credit: https://neurohive.io/en/popular-networks/vgg16/.

128 nodes was added with the ReLU activation function and, finally, a five-node output layer with softmax activation was used [68].

*2) Convolutional neural network based on the xception architecture:* Another architecture assessed in this work is Xception. Xception was developed at Google by F. Chollet. It is based on the concept of "inception modules," and is endowed with 36 convolutional layers for feature extraction [69]. This architecture seems very promising, as it outperformed other precedent networks with similar number of parameters on large data sets such as ImageNet [56] and another Google internal data set denoted by JFT.

Like VGG16, the Xception architecture can be also implemented through the Keras framework [70]. In its default implementation, the size of the input images is 299 pixels × 299 pixels, with 3 channels. However, disabling the default top layer of this CNN allows for the use of images with other sizes. We opted to disable the default top layer, in order to keep the same input size of 300 pixels × 300 pixels, and replaced the disabled layer by another fully connected layer with 128 nodes.

In the resulting architecture, the image is processed through an array of 14 blocks of successive convolutional layers. Each block is composed differently, combining *depth-wise separable* convolutional layers with ReLU activation functions, and $3\times3$ max-pooling layers. As with the previous architecture (Section II-A1), the output of the convolutional layers was passed through a filter of the average pool type, with a size of $2 \times 2$ pixels and a fully connected layer of 128 nodes with the ReLU activation function. Likewise, the output layer had five nodes with softmax activation.

*3) Training and testing sets:* In this experiment, the whole data set consisted of 1500 images, which included 300 images of type UI, 300 images of type DI, 300 images of type LI, 300 images of type RI, and 300 images of type N. The training set was constructed including 250 images of each class, while the remaining images (50 from each class) constituted the testing set.

In order to generate the images of class N (null inclination), the mobile was held parallel to the surface, with an allowed error up to 3 deg (i.e., $i \leq 3$ deg). On the other hand, the images that were not intended to be of class N (i.e., with inclination) were taken with $i$ within a range of 30–70 deg. With respect to $r$, for every class other than N, we allowed it to include values of $r$ centered at their nominal values $r_0$ (0, 90, 180, or, 270 deg) and within intervals of $r \in [r_0 \pm 5\,\text{deg}]$.

*4) Training process:* For both CNN architectures, the weights that were optimized were those of the fully connected layers, while the weights of the convolutional layers were set as the default pre-trained values.

The optimization of the weights was achieved using the Adam optimizer [71], for which the learning rate $lr$ was scheduled as $lr(k) = 0.001/(1 + d \cdot k)$, where $k$ is the iteration number and $d = 0.00002857142$. The loss function $\mathcal{J}$ was categorical cross-entropy. Both models were trained for 40 epochs with batches of 30 observations.

### B. Experiment II

Experiment II was considered as a natural step forward from Experiment I. In this case, more refined categories for $i$ were pursued. We intended to obtain an algorithm that could distinguish incidence angles in three classes: $i_0 = 0$ deg, $i_0 = 20$ deg, and $i_0 = 40$ deg. The sub-index 0 indicates the nominal value for each corresponding class. The real images were taken with some errors allowed (for $i$ up to 3 deg) from their corresponding nominal values. In other words, the class with $i_0 = 0$ deg actually contained angles $i \leq 3$ deg. The class

of $i_0 = 20$ deg implied $17 \deg \leq i \leq 23 \deg$, while that of $i_0 = 40$ deg included angles in the range $37 \deg \leq i \leq 43 \deg$. For each class with $i_0 > 0$ deg, the algorithm had to distinguish between $r_0 = 0$ deg, $r_0 = 90$ deg, $r_0 = 180$ deg, and $r_0 = 270$ deg. In total, there were nine classes, denoted by the following: 'N' ($i_0 = 0$ deg), '20U' ($i_0 = 20$ deg and $r_0 = 0$ deg), '20R' ($i_0 = 20$ deg and $r_0 = 90$ deg), '20D' ($i_0 = 20$ deg and $r_0 = 180$ deg), '20L' ($i_0 = 20$ deg and $r_0 = 270$ deg), '40U' ($i_0 = 40$ deg and $r_0 = 0$ deg), '40R' ($i_0 = 40$ deg and $r_0 = 90$ deg), '40D' ($i_0 = 40$ deg and $r_0 = 180$ deg), and '40L' ($i_0 = 40$ deg and $r_0 = 270$ deg).

In this experiment, the architectures described in Section II-A1 and Section II-A2 were initially attempted, but their results were not promising. Hence, a third architecture was pursued. This architecture was proposed in [68] and, due to its similarities to the VGGNet architecture [62], it is referred to as the compact version of VGGNet, named *SmallerVGGNet*. The following section provides a description of it.

In this experiment, rather than considering nine mutually exclusive classes, the problem was addressed as a multi-label classification process. Each image was assigned to two labels. One label indicated the incidence angle: either $i_0 = 0$ deg, $i_0 = 20$ deg, or $i_0 = 40$ deg. For the cases with $i_0 > 0$ deg, another label expressing the angle $r_0$ was associated with them: either $r_0 = 0$ deg, $r_0 = 90$ deg, $r_0 = 180$ deg, or $r_0 = 270$ deg. Thereby, by using the CNN to classify the images on $i_0$ and $r_0$, all of the cases of interest were addressed.

*1) SmallerVGGNet architecture:* A schematic representation of this architecture can be found in [68]. Following the input layer, it has a first convolutional layer with 32 kernels of size $3 \times 3$, activated by the Rectified Linear Unit (ReLU) function. It uses a padding of 1 and stride of 1. This layer is followed by a MaxPool layer of size $3 \times 3$ and a stride of $3 \times 3$. A dropout scheme with a rate of 25% is applied before the next convolutional layer. Then, there are two convolutional layers, each of which has 64 filters of size $3 \times 3$, and a ReLU activation function. These are followed by another MaxPool layer of size and stride $2 \times 2$. Another dropout layer with a rate of 25% was applied before the next convolutional layer. These layers are followed by another two convolutional layers, with 128 kernels each, where each kernel had a size of $3 \times 3$. These convolutional layers are activated with ReLU functions, and are concatenated to another MaxPool layer, of size and stride $2 \times 2$, and a dropout layer with rate 25%.

Following the aforementioned layers, there is a fully connected layer of 1024 nodes with ReLU activation, a dropout layer with a rate of 50%, and the final output layer activated with sigmoid functions. It is important to note that, for Experiment II, the problem was tackled as a multi-label classification one, categorizing both angles $i$ and $r$ independently.

*2) Training and testing sets:* In this experiment, the training set consisted of:

- 1252 images of type N;
- 293 images of type 20U;
- 282 images of type 20R;
- 289 images of type 20D;
- 292 images of type 20L;
- 289 images of type 40U;
- 301 images of type 40R;
- 284 images of type 40D; and
- 288 images of type 40L.

Meanwhile, the testing set contained:

- 220 images of type N;
- 45 images of type 20U;
- 56 images of type 20R;
- 49 images of type 20D;
- 46 images of type 20L;
- 55 images of type 40U;
- 43 images of type 40R;
- 60 images of type 40D; and
- 56 images of type 40L.

The number of images of type N might seem much higher than that for the other classes. This is due to the fact that, when the initial set of images was augmented by rotating them, all of the rotated images within the category N remained in the same category. All these images were used with the intention to provide images in category N with different roll angles.

*3) Training process:* The CNN was trained for initially 30 epochs, with batches of 32 images each. Optimization of the weights was achieved by using the built-in Adam optimizer [71], for which the learning rate $lr$ was scheduled as $lr(k) = 0.001 / (1 + d \cdot k)$, where $k$ is the iteration number and $d = 0.00002857142$. After the first 30 epochs, the CNN was retrained for another 10 epochs, but with a learning rate given by $lr(k) = 0.0005 / (1 + d \cdot k)$.

### C. Images and Surfaces Considered

The images used in this work for training show multiple type of surfaces with different textures. They represent ground terrains where a lander might accomplish touchdown. The appearance of surfaces where a space probe might land can be very diverse. There might be areas with multiple protruding rocks, or areas that are mostly plain [72]. For VTOL aircrafts, the landing surface could be also very varied, including grass, concrete, and paving stones. We intended to obtain an algorithm that can work with multiple textures, including some that can be found on the Earth's ground and others that are not of ground-like types. Including images with multiple textures allowed the CNN to learn features that can provide indications about different inclinations, even across images of different textures. Samples of these surfaces are shown in Fig. 6c, and Fig. 9a–Fig. 9l.

The images were taken using a mobile telephone camera. The device was a Samsung® S6 Edge. The resolution of the camera was 16 megapixels. While taking the images, the mobile was held manually, at a distance between 30 and 100 cm from the surface. However, for missions, the distance

(a) Grass surface, $r = 270$ deg.



(b) Concrete surface, $r = 90$ deg.



(c) Wood surface, $r = 0$ deg.



(d) Tiles surface, $r = 180$ deg.



(e) Wall cover, $r = 90$ deg.



(f) Soiled concrete surface with rocks, $r = 0$ deg.



(g) Tiles surface (other type), with $r = 0$ deg.



(h) Tiles surface (other type), $r = 0$ deg.



(i) Table cover, $r = 270$ deg.



(j) Wood floor, $r = 0$ deg.



(k) Granite surface, $r = 270$ deg.



(l) Table cover, $r = 270$ deg.

Fig. 9. Different angles $r$ for images with $i > 0$ deg.

from the ground at which images could be taken may be highly diverse. Two important factors that would dictate the appropriate distances are the resolution of the sensor of the camera and its lenses, which entail the size of each portion of ground represented by each pixel in the image. Hence, in this work, the distance from the ground at which the images were taken were arbitrarily set, as the main goal was to demonstrate the methodology, rather than obtaining a production-level algorithm for a specific mission. Furthermore, the incidence angle does not depend on the distance to the ground at which the images are taken. Therefore, we expect that the features that characterize the angles $i$ and $r$ of an image could be learned by a CNN across multiple images from different distances.

*1) Image Pre-processing and augmentation:* In general, convolutional neural networks require inputs of a pre-determined size. The size $h = w = 300$ pixels of the input of the CNN was arbitrarily chosen. These values were considered to provide images that were as large as possible (in order to provide the training process with as much information as possible), but still allowing for a training process that could be carried out entirely without being prematurely terminated due to a RAM memory shortage.

Once the images were resized to a common size, every image was rotated clockwise three times (i.e., by 90, 180, and 270 deg), in order to augment the data set. Moreover, as the original images were of RGB type (i.e., their color channels were in the order of red, green, and blue), the set was augmented by converting them to GBR (i.e., green, blue, and red, in that order).

Since the images were RGB (three color channels), each image was represented by a tensor of dimension $300 \times 300 \times 3$. Once the final set of resized and rotated images was completely defined, the intensity value corresponding to each pixel, for each of the color channels, was divided by 255, in order to scale their values into the range $[0, 1]$.

### III. Results and Discussion

#### A. Experiment I

In this experiment, both CNN architectures were trained and tested with the same training and testing sets. Recalling the classification categories for Experiment I, Table I and Table II show the confusion matrices obtained with the testing set, for each of the two architectures. In these tables, U indicates an upward inclination, L denotes leftward inclination, D is downward inclination, and R represents rightward inclination.

TABLE I. Confusion Matrix for VGG16 CNN

| | | Predicted Classes | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | U | L | D | R | N |
| | U | 43 | 0 | 0 | 3 | 4 |
| | L | 3 | 40 | 0 | 4 | 3 |
| True Labels | D | 1 | 2 | 41 | 1 | 5 |
| | R | 2 | 3 | 2 | 39 | 4 |
| | N | 1 | 1 | 2 | 6 | 40 |

A few observations can be drawn from Table I and Table II. First, the number of true positive cases for each category, in both tables, strongly suggests that it is feasible to train an algorithm to classify the images, according to the categories defined in this work. This means that, from individual monocular images, CNNs can extract useful information to determine whether the terrain under a landing vehicle is inclined with

TABLE II. Confusion Matrix for Xception CNN

| | | Predicted Classes | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | U | L | D | R | N |
| | U | 43 | 1 | 0 | 1 | 5 |
| | L | 2 | 41 | 1 | 3 | 3 |
| True Labels | D | 1 | 0 | 41 | 0 | 8 |
| | R | 0 | 4 | 8 | 34 | 4 |
| | N | 0 | 5 | 3 | 2 | 40 |

TABLE III. Precision and Recall for the VGG16-based CNN

| | Precision | Recall | $F_1$ Score |
| --- | --- | --- | --- |
| U | 0.86 | 0.86 | 0.86 |
| L | 0.87 | 0.80 | 0.83 |
| D | 0.91 | 0.82 | 0.86 |
| R | 0.74 | 0.78 | 0.76 |
| N | 0.71 | 0.80 | 0.75 |

TABLE IV. Precision and Recall for the Xception-based CNN

| | Precision | Recall | $F_1$ Score |
| --- | --- | --- | --- |
| U | 0.93 | 0.86 | 0.90 |
| L | 0.80 | 0.82 | 0.81 |
| D | 0.77 | 0.82 | 0.80 |
| R | 0.85 | 0.68 | 0.76 |
| N | 0.67 | 0.80 | 0.73 |

respect to the attitude of the vehicle, as well as the direction of the inclination.

We found that the CNN based on VGG16 performed slightly better than its Xception counterpart. Although the numbers were somewhat similar in magnitude, the 39 correctly predicted cases of rightward inclination against the 34, indicated the advantage of VGG16 over Xception. For statistical comparison, we observed the metrics of precision, recall, and $F_1$ score, as they constitute natural manners to evaluate the performance of classification algorithms. Considering the

Fig. 10. Evolution of the training of the VGG16-based CNN.

TABLE V. CONFUSION MATRIX: EXPERIMENT II

| | | | | | Predicted Classes | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | N | 20U | 20R | 20D | 20L | 40U | 40R | 40D | 40L |
| | **N** | 220 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | **20U** | 7 | 34 | 1 | 1 | 0 | 2 | 0 | 0 | 0 |
| | **20R** | 4 | 0 | 49 | 3 | 0 | 0 | 0 | 0 | 0 |
| | **20D** | 9 | 0 | 0 | 40 | 0 | 0 | 0 | 0 | 0 |
| **True Labels** | **20L** | 4 | 4 | 0 | 0 | 38 | 0 | 0 | 0 | 0 |
| | **40U** | 5 | 0 | 0 | 0 | 0 | 50 | 0 | 0 | 0 |
| | **40R** | 4 | 0 | 1 | 0 | 0 | 0 | 38 | 0 | 0 |
| | **40D** | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 54 | 0 |
| | **40L** | 8 | 0 | 0 | 0 | 4 | 0 | 0 | 1 | 43 |

precision, recall, and $F_1$ scores displayed in Table III and Table IV, none of the architectures outperformed the other in every other metric.

Fig. 10 illustrates the progress in the prediction capability

TABLE VI. PRECISION AND RECALL FOR EXPERIMENT II

|  | Precision | Recall | $F_1$ Score |
|---|---|---|---|
| **N** | 0.82 | 0.76 | 0.90 |
| **20U** | 0.89 | 0.80 | 0.82 |
| **20R** | 0.96 | 0.88 | 0.92 |
| **20D** | 0.91 | 0.82 | 0.86 |
| **20L** | 0.90 | 0.83 | 0.86 |
| **40U** | 0.96 | 0.91 | 0.93 |
| **40R** | 1 | 0.88 | 0.94 |
| **40D** | 0.98 | 0.90 | 0.94 |
| **40L** | 1 | 0.77 | 0.87 |

of the CNN as the training advanced. It was observed that, after 25 epochs of training, the loss function of the validation set stopped decreasing. We also observed that the loss function computed over the training set and that over the validation set diverged after approximately 10 epochs, which might be indicative of overfitting. This could be resolved by adding more images to the training set.

### B. Experiment II

For the performance obtained in Experiment II, using the architecture described in Section II-B1, Table V shows the distribution of classifications obtained for the testing set, while Table VI indicates the precision, recall, and $F_1$ score statistics obtained. These numbers suggest that the algorithm, indeed, learned to distinguish the classes to which each image belonged. This supports our work towards the next step, which is generating a classifier that can provide estimates of the angles $i$ and $r$, but within many more classes; thus, describing the measured angles more precisely.

## IV. CONCLUSION

In this work, we explored the feasibility of using convolutional neural networks to evaluate the slope of the terrain under a lander, when it is about to touch ground. This capability may be essential for certain missions, where autonomous landers must accomplish landing maneuvers in unknown terrains.

Two experiments were described, where Experiment II was considered as a natural extension of Experiment I. The latter demonstrated the feasibility of using CNNs to classify image angles into five categories, including normal and four categories with $i > 0$. The former exhibited that it is also possible to quantify the angles $i$ into more than binary categories. The next step will be to train algorithms that can classify the images into many more categories, or treat the problem as a regression one, in which the outputs are real numbers for the angles $r$ and $i$.

It is important to mention that this work constitutes the first step in a wider project, whose ultimate goal is developing algorithms to precisely estimate angles between landers and the terrain underneath them from monocular images.

## REFERENCES

[1] Kulumani, S.; Takami, K.; Lee, T. Geometric control for autonomous landing on asteroid Itokawa using visual localization. Proceedings of the AAS/AIAA Astrodynamics Specialist Conference, Stevenson, Washington, USA, 20–24 August 2017; Univelt, Inc.: Escondido, CA, USA, 2017.

[2] Gaudet, B.; Linares, R.; Furfaro R. Terminal adaptive guidance via reinforcement meta-learning: Applications to autonomous asteroid close-proximity operations. *Acta Astronautica* 2020, *171*, 1–13. https://doi.org/10.1016/j.actaastro.2020.02.036.

[3] Stacey, N.; D'Amico S. Autonomous swarming for simultaneous navigation and asteroid characterization. Proceedings of the AAS/AIAA Astrodynamics Specialist Conference, Snowbird, UT, USA, 19–23 August 2018; Univelt, Inc.: Escondido, CA, USA, 2018.

[4] Sandino, J.; Vanegas, F.; Maire, F.; Caccetta, P.; Sanderson, C.; Gonzalez, F. UAV framework for autonomous onboard navigation and people/object detection in cluttered indoor environments. *Remote Sens* 2020, *12*, 3386. https://doi.org/10.3390/rs12203386.

[5] Albattah, W.; Masood, M.; Javed, A.; et al. Custom CornerNet: A drone-based improved deep learning technique for large-scale multiclass pest localization and classification. *Complex Intell Syst* 2023, *9*, 1299--1316. https://doi.org/10.1007/s40747-022-00847-x.

[6] Bhaskaran, S.; Nandi, S.; Broschart, S.; Wallace, M.; Cangahuala, L.A.; Olson, C. Small body landing accuracy using in-situ navigation. 2011. Available online: https://trs.jpl.nasa.gov/bitstream/handle/2014/41886/11-0341.pdf?sequence=1&isAllowed=y (accessed on 01 April 2021).

[7] Silva, M.F.; Cerqueira, A.S.; Vidal, V.F.; Honório, L.M.; Santos, M.F.; Oliveira, E.J. Landing area recognition by image applied to an autonomous control landing of VTOL aircraft. Proceedings of the International Carpathian Control Conference (ICCC), Sinaia, Romania, 28-31 May 2017; IEEE: New York, NY, USA, 2017.

[8] Aziz, M.Z.; Mertsching, B. Survivor search with autonomous UGVs using multimodal overt attention. Proceedings of the Safety Security and Rescue Robotics, Bremen, Germany, 26-30 July 2010; IEEE: New York, NY, USA, 2011.

[9] Kawano, I.; Mokuno, M.; Kasai, T.; Suzuki, T. First autonomous rendezvous using relative GPS navigation by ETS-VII. *Navig* 2001, *48*, 49–56. https://doi.org/10.1002/j.2161-4296.2001.tb00227.x.

[10] Furfaro, R.; Cersosimo, D.; Wibben, D.R. Asteroid precision landing via multiple sliding surfaces guidance techniques. *J of Guid, Control, and Dyn* 2013, *36*, 1075–1092. https://doi.org/10.2514/1.58246.

[11] Tognon, M.; Testa, A.; Rossi, E.; Franchi, A. Takeoff and landing on slopes via inclined hovering with a tethered aerial robot. Proceedings of the International Conference on Intelligent Robots and Systems (IROS), Daejeon, South Korea, 09-14 October 2016; IEEE: New York, NY, USA, 2016.

[12] Yano, H.; Kubota, T.; Miyamoto, H.; et al. Touchdown of the Hayabusa spacecraft at the Muses Sea on Itokawa. *Sci* 2006, *312*, 1350–1353. https://doi.org/10.1126/science.1126164.

[13] Liebe, C.C.; Abramovici, A.; Bartman, R.K.; et al. Laser radar for spacecraft guidance applications. Proceedings of the Aerospace Conference (Cat. No.03TH8652), Big Sky, MT, USA, 8-15 March 2003; IEEE: New York, NY, USA, 2003.

[14] Shimkin, P.E.; Baskakov, A.I.; Komarov, A.A.; Ka, M. Safe helicopter landing on unprepared terrain using onboard interferometric radar. *Sensors* 2020, *20*, 2422. https://doi.org/10.3390/s20082422.

[15] Ansar, A.; Cheng, Y. Vision technologies for small body proximity operations. Proceedings of the International Symposium on Artificial Intelligence, Robotics and Automation in Space (i-SAIRAS), Sapporo, Japan, 29 August - 1 September 2010; European Space Agency (ESA): Paris, France, 2010.

[16] Segal, S.; Carmi, A.; Gurfil, P. Stereovision-based estimation of relative dynamics between noncooperative satellites: Theory and experiments. *IEEE Trans on Control Syst Technol* 2014, *22*, 568–584. https://doi.org/10.1109/TCST.2013.2255288.

[17] Sharma, S.; Beierle, C.; D'Amico, S. Pose estimation for non-cooperative spacecraft rendezvous using convolutional neural networks. Proceedings of the Aerospace Conference, Big Sky, MT, USA, 3-10 March 2018; IEEE: New York, NY, USA, 2018.

[18] Herissé, B.; Hamel, T.; Mahony, R.; Russotto, F.X. The landing problem of a VTOL unmanned aerial vehicle on a moving platform using optical flow. Proceedings of the International Conference on Intelligent Robots and Systems, Taipei, Taiwan, USA, 18-22 October 2010; IEEE: New York, NY, USA, 2010.

[19] Herissé, B.; Hamel, T.; Mahony, R.; Russotto, F.X. Landing a VTOL unmanned aerial vehicle on a moving platform using optical flow. *IEEE Trans on Robotics* 2012, *28*, 77–89. https://doi.org/10.1109/TRO.2011.2163435.

[20] Rosa, L.; Hamel, T.; Mahony, R.; Samson, C. Optical-flow based strategies for landing VTOL UAVs in cluttered environments. Proceedings of the World Congress of the International Federation of Automatic Control (IFAC), Cape Town, South Africa, 24-29 August 2014; Elsevier: Amsterdam, Netherlands, 2016.

[21] Green, W.E.; Oh, P.Y. Optic-flow-based collision avoidance. *IEEE Robotics Autom Mag* 2008, *15*, 96–103. https://doi.org/10.1109/MRA.2008.919023.

[22] Beyeler, A.; Zufferey, J.C.; Floreano, D. Vision-based control of near-obstacle flight. *Auton Robots* 2009, *27*, 201–219. https://doi.org/10.1007/s10514-009-9139-6.

[23] Padhy, R.P.; Verma, S.; Ahmad, S.; Choudhury, S.K.; Sa, P.K. Deep neural network for autonomous UAV navigation in indoor corridor environments. *Procedia Comput Sci* 2018, *133*, 643–650. https://doi.org/10.1016/j.procs.2018.07.099.

[24] Lu, Y.; Xue, Z.; Xia, G.S.; Zhang, L. A survey on vision-based UAV navigation. *Geospat Inf Sci* 2018, *21*, 21–32. https://doi.org/10.1080/10095020.2017.1420509.

[25] Mohan, K.; M, P.; S, S.; et al. LIDAR based landing site identification and safety estimation for inter planetary missions. Proceedings of the International Conference on Control, Communication and Computing (ICCC), Thiruvananthapuram, India, 19-21 May 2023; IEEE: New York, NY, USA, 2023.

[26] Xie, H.; Tang, H.; Jin, Y.; et al. An improved surface slope estimation model using space-borne laser altimetric waveform data over the Antarctic ice sheet. *IEEE Geoscience and Remote Sensing Letters* 2022, *19*, 1–5. 10.1109/LGRS.2021.3124224.

[27] Arai, K. Ground control point generation from simulated SAR image derived from digital terrain model and its application to texture feature extraction. *International Journal of Advanced Computer Science and Applications* 2021, *12*. http://dx.doi.org/10.14569/IJACSA.2021.0120112.

[28] Lee, H.Y.; Kim, T.; Park, W.; Lee, H.K. Extraction of digital elevation models from satellite stereo images through stereo matching based on epipolarity and scene geometry. *Image and Vis Comput* 2003, *21*, 789–796. https://doi.org/10.1016/S0262-8856(03)00092-1.

[29] Ajayi, O.G.; Salubi, A.A.; Angbas, A.F.; Odigure, M.G. Generation of accurate digital elevation models from UAV acquired low percentage overlapping images. *Int J of Remote Sens* 2017, *38*, 3113–3134. https://doi.org/10.1080/01431161.2017.1285085.

[30] Krupnik, A. Accuracy assessment of automatically derived digital elevation models from SPOT images. *Photogramm Eng and Remote Sens* 2000, *66*, 1017–1023.

[31] Zhu, X.; Nie, S.; Wang, C.; Xi, X.; Li, D.; Li, G.; Wang, P.; Cao, D.; Yang, X. Estimating terrain slope from ICESat-2 data in forest environments. *Remote Sens* 2020, *12*, 3300. https://doi.org/10.3390/rs12203300.

[32] Stevens, K.A. The information content of texture gradients. *Biological Cybern* 1981, *42*, 95–105. https://doi.org/10.1007/BF00336727.

[33] Ikeuchi, K. Shape from regular patterns. *Artif Intell* 1984, *22*, 49–75. https://doi.org/10.1016/0004-3702(84)90025-0.

[34] Knill, D.C. Surface orientation from texture: Ideal observers, generic observers and the information content of texture cues. *Vis Res* 1998, *38*, 1655–1682. https://doi.org/10.1016/S0042-6989(97)00324-6.

[35] Kanatani, K.; Chou, T.C. Shape from texture: General principle. *Artif Intell* 1989, *38*, 1–48. https://doi.org/10.1016/0004-3702(89)90066-0.

[36] Malik, J.; Rosenholtz, R. Recovering surface curvature and orientation from texture distortion: A least squares algorithm and sensitivity analysis. Proceedings of the European Conference on Computer Vision (ECCV), Stockholm, Sweden, 2-6 May 1994; Springer: Berlin, Germany, 2005.

[37] Malik, J.; Rosenholtz, R. Computing local surface orientation and shape from texture for curved surfaces. *Int J of Comput Vis (IJCV)* 1997, *23*, 149–168. https://doi.org/10.1023/A:1007958829620.

[38] Prince, M.; Alsuhibany, S. A.; Siddiqi, N. A. A step towards the optimal estimation of image orientation. *IEEE Access* 2019, *7*, 185750–185759. https://doi.org/10.1109/ACCESS.2019.2959666.

[39] Romdhani, S.; Vetter, T. Estimating 3D shape and texture using pixel intensity, edges, specular highlights, texture constraints and a prior. Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR), San Diego, CA, USA, 20-25 June 2005; IEEE: New York, NY, USA, 2005.

[40] Saxena, A.; Schulte, J.; Ng, A.Y. Depth estimation using monocular and stereo cues. Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI), Hyderabad, India, 6-12 January 2007; IJCAI: CA, USA, 2007.

[41] Saxena, A.; Chung, S.H.; Ng, A.Y. 3-D depth reconstruction from a single still image. *Int J of Comput Vis (IJCV)* 2008, *76*, 53–69. https://doi.org/10.1007/s11263-007-0071-y.

[42] Handa, A.; Sharma, P. Real-time depth estimation from a monocular moving camera. Proceedings of the International Conference on Contemporary Computing (IC3), Noida, India, 6-8 August 2012; Springer: Berlin, Germany, 2012.

[43] Srikakulapu, V.; Kumar, H.; Gupta, S.; Venkatesh, K.S. Depth estimation from single image using defocus and texture cues. Proceedings of the National Conference on Computer Vision, Pattern Recognition, Image Processing and Graphics (NCVPRIPG), Patna, India, 16-19 December 2015; IEEE: New York, NY, USA, 2016.

[44] Lee, J.; Kim, C. Monocular depth estimation using relative depth maps. Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, 15-20 June 2019; IEEE: New York, NY, USA, 2020.

[45] Mertan, A.; Duff, D.J.; Unal, G. Relative depth estimation as a ranking problem. *arXiv* 2020, arXiv:2010.06944.

[46] Meingast, M.; Geyer, C.; Sastry, S. Vision based terrain recovery for landing unmanned aerial vehicles. Proceedings of the Conference on Decision and Control (CDC) (IEEE Cat. No.04CH37601), Nassau, Bahamas, 14-17 December 2004; IEEE: New York, NY, USA, 2005.

[47] Rothrock, B.; Kennedy, R.; Cunningham, C.; Papon, J.; Heverly, M.; Ono, M. SPOC: Deep learning-based terrain classification for Mars rover missions. Proceedings of the AIAA Space, Long Beach, CA, USA, 13-16 September 2016; AIAA: Reston, VA, USA, 2016.

[48] Bellone, M.; Messina, A.; Reina, G. A new approach for terrain analysis in mobile robot applications. Proceedings of the International Conference on Mechatronics (ICM), Vicenza, Italy, 27 February - 1 March 2013; IEEE: New York, NY, USA, 2013.

[49] Gonzalez, R.; Iagnemma, K. DeepTerramechanics: Terrain classification and slip estimation for ground robots via deep learning. *arXiv* 2018, arXiv:1806.07379.

[50] Valada, A.; Burgard, W. Deep spatiotemporal models for robust proprioceptive terrain classification. *Int J of Robotics Res* 2017, *36*, 1521–1539. https://doi.org/10.1177/0278364917727062.

[51] Dreiseitl, S.; Ohno-Machado, L. Logistic regression and artificial neural network classification models: A methodology review. *J of Biomed Informatics* 2002, *35*, 352–359. https://doi.org/10.1016/S1532-0464(03)00034-0.

[52] Chaudhuri, B.B.; Bhattacharya, U. Efficient training and improved

performance of multilayer perceptron in pattern classification. *Neurocomputing* 2000, *34*, 11–27. https://doi.org/10.1016/S0925-2312(00)00305-2.

[53]   Orhan, U.; Hekim, M.; Ozer, M. EEG signals classification using the K-means clustering and a multilayer perceptron neural network model. *Expert Syst with Appl* 2011, *38*, 13475–13481. https://doi.org/10.1016/j.eswa.2011.04.149.

[54]   Al-Saffar, A.A.M.; Tao, H.; Talab, M.A. Review of deep convolution neural network in image classification. Proceedings of the International Conference on Radar, Antenna, Microwave, Electronics, and Telecommunications (ICRAMET), Jakarta, Indonesia, 23-24 October 2017; IEEE: New York, NY, USA, 2018.

[55]   Khan, S.; Rahmani, H.; Shah, S.A.A.; Bennamoun, M. *A Guide to Convolutional Neural Networks for Computer Vision*, 1st ed.; Morgan & Claypool: San Rafael, CA, USA, 2018.

[56]   Russakovsky, O.; Deng, J.; Su, H.; et al. ImageNet large scale visual recognition challenge. *Int J Comput Vis (IJCV)* 2015, *115*, 211–252. https://doi.org/10.1007/s11263-015-0816-y.

[57]   Chen, Y.; Xie, H.; Shin, H. Multi-layer fusion techniques using a CNN for multispectral pedestrian detection. *IET Comput Vis* 2018, *12*, 1179–1187. https://doi.org/10.1049/iet-cvi.2018.5315.

[58]   Zhang, S.; Wen, L.; Bian, X.; Lei, Z.; Li, S.Z. Occlusion-aware R-CNN: Detecting pedestrians in a crowd. Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, 8-14 September 2018; Springer: Berlin, Germany, 2018.

[59]   Goel, T.; Murugan, R.; Mirjalili, S.; Chakrabartty, D.K. OptCoNet: An optimized convolutional neural network for an automatic diagnosis of COVID-19. *Appl Intell* 2021, *51*, 1351–1366. https://doi.org/10.1007/s10489-020-01904-z.

[60]   Misra, S.; Jeon, S.; Lee, S.; Managuli, R.; Jang, I.; Kim, C. Multi-channel transfer learning of chest X-ray images for screening of COVID-19. *Electron* 2020, *9*, 1388. https://doi.org/10.3390/electronics9091388.

[61]   Jain, R.; Gupta, M.; Taneja, S.; Hemanth, D.J. Deep learning based detection and analysis of COVID-19 on chest X-ray images. *Appl Intell* 2021, *51*, 1690–1700. https://doi.org/10.1007/s10489-020-01902-1.

[62]   Simonyan, K.; Zisserman A. Very deep convolutional networks for large-scale image recognition. Proceedings of the International Conference on Learning Representations (ICLR), San Diego, CA, USA, 7-9 May 2015; arXiv: New York, NY, USA, 2015.

[63]   Tindall, L.; Luong, C.; Saad, A. Plankton classification using VGG16 network. 2017. Available online: https://pdfs.semanticscholar.org/7cb1/a0d0d30b4567b771ad7ae265ab0e935bc41c.pdf (accessed on 01 April 2021).

[64]   Ashraf, K.; Wu, B.; Iandola, F.; Moskewicz, M.; Keutzer, K. Shallow networks for high-accuracy road object-detection. *arXiv* 2016, arXiv:1606.01561.

[65]   Antony, J.; McGuinness, K.; O'Connor, N.; Moran, K. Quantifying radiographic knee osteoarthritis severity using deep convolutional neural networks. Proceedings of the International Conference on Pattern Recognition (ICPR), Cancun, Mexico, 4-8 December 2016; IEEE: New York, NY, USA, 2016.

[66]   Keras. VGG16 and VGG19. 2015. Available online: https://keras.io/api/applications/vgg/#vgg16-function (accessed on 01 April 2021).

[67]   TensorFlow. Available online: https://www.tensorflow.org/ (accessed on 01 April 2021).

[68]   PyImageSearch. Available online: https://www.pyimagesearch.com (accessed on 01 April 2021).

[69]   Chollet, F. Xception: Deep learning with depthwise separable convolutions. Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR), Honolulu, HI, USA, 21-26 July 2017; IEEE: New York, NY, USA, 2017.

[70]   Keras. Xception. 2017. Available online: https://keras.io/api/applications/xception/ (accessed on 01 April 2021).

[71]   Kingma, D.; Ba, J. Adam: A method for stochastic optimization. Proceedings of the International Conference on Learning Representations (ICLR), San Diego, CA, USA, 7-9 May 2015; arXiv: New York, NY, USA, 2015.

[72]   Saito, J.; Miyamoto, H.; Nakamura, R.; et al. Detailed images of asteroid 25143 Itokawa from Hayabusa. *Sci* 2006, *312*, 1341–1344. https://doi.org/10.1126/science.1125722.

# Transformative Automation: AI in Scientific Literature Reviews

Kirtirajsinh Zala[1], Biswaranjan Acharya[2], Madhav Mashru[3],
Damodharan Palaniappan[4], Vassilis C. Gerogiannis[5], Andreas Kanavos[6], Ioannis Karamitsos[7]
Department of Information Technology, Marwadi University, Rajkot, Gujarat 360003, India[1,4]
Department of Computer Engineering -AI & BDA, Marwadi University, Rajkot, Gujarat 360003, India[2]
Faculty of Engineering, Marwadi Education Foundation's Group of Institutions, Rajkot, Gujarat 360003, India[3]
Department of Digital Systems, University of Thessaly, Larissa, Greece[5]
Department of Informatics, Ionian University, Corfu, Greece[6]
Research and Graduate Department, Rochester Institute of Technology, Dubai, UAE[7]

*Abstract*—This paper investigates the integration of Artificial Intelligence (AI) into systematic literature reviews (SLRs), aiming to address the challenges associated with the manual review process. SLRs, a crucial aspect of scholarly research, often prove time-consuming and prone to errors. In response, this work explores the application of AI techniques, including Natural Language Processing (NLP), machine learning, data mining, and text analytics, to automate various stages of the SLR process. Specifically, we focus on paper identification, information extraction, and data synthesis. The study delves into the roles of NLP and machine learning algorithms in automating the identification of relevant papers based on defined criteria. Researchers now have access to a diverse set of AI-based tools and platforms designed to streamline SLRs, offering automated search, retrieval, text mining, and analysis of relevant publications. The dynamic field of AI-driven SLR automation continues to evolve, with ongoing exploration of new techniques and enhancements to existing algorithms. This shift from manual efforts to automation not only enhances the efficiency and effectiveness of SLRs but also marks a significant advancement in the broader research process.

*Keywords*—*Artificial intelligence; systematic literature review; scholarly data analysis; machine learning algorithms; natural language processing; scientific publication automation*

## I. INTRODUCTION

Artificial intelligence (AI) has emerged to alleviate humans from repetitive tasks that demand specific human skills. Like any other field, scientific endeavors benefit from powerful algorithms to expedite and enhance outcomes. Initiating a new research project typically involves a thorough investigation of relevant scholarly publications to comprehend the landscape and identify activities significant for addressing similar or related issues. The process of gathering documents, when performed without prior training or well-defined parameters, may lead to the omission of significant contributions [29]. A comprehensive approach to searching and analyzing literature can help reduce the likelihood of bias and inaccuracy in research [24], [28].

A systematic literature review (SLR) is a secondary investigation that assesses existing research, employing a widely recognized procedure to identify related articles, extract pertinent details, and present their main findings in an organized manner [33]. It is anticipated that a published literature review will deliver a comprehensive summary of a corresponding research subject, often providing a historical perspective that facilitates the identification of research trends and unresolved issues. Literature reviews are now a fundamental component of many scientific fields, including medicine (with 13,510 published reviews) and computer science (with 6,342) [47].

Conducting a literature review is known to be time-consuming, especially when addressing a vast research subject. In recent years, various systematic literature review (SLR)-related tools have been developed for diverse purposes [47]. These tools can automate digital database searches, designate relevant outcomes based on inclusion criteria, and provide visual support for analyzing information from works' authors and their citations, among other capabilities. Particularly, the automation of the SLR process is gaining attention in the field of computer science research, offering strategies to construct search phrases and retrieve publications semi-automatically or manually from relevant scientific databases [76]. The utilization of automated methods has proven to save time and costs in selecting relevant articles [11], or providing a summary of the findings [71]. However, some authors argue that the usefulness of these automated tools is limited by their steep learning curve and the lack of research analyzing the advantages they offer [74].

This paper focuses on the computerized and automated operation of SLR tasks, replacing manual labor with ML as the primary driver. The goal is to enhance the capability of automated review processes and technologies with some additional understanding and suggestions. The initial application of AI methods to automate SLR tasks occurred in 2006 [12], where it was suggested that neural networks could be used to automate the selection of relevant articles. Initial resistance to this idea stemmed from concerns regarding the use of data gleaned from secondary sources through text mining [51].

Following this concept, previous works by other researchers have delved into powerful text mining techniques [52], [58], [65]. Recent innovations in the field include the integration of ML and natural language processing (NLP) techniques [27], [76]. Considering the repetitive tasks involved in a SLR methodology, the capabilities of AI for analyzing scientific literature are vast. However, it's crucial not to devalue the role of human involvement in this process, as humans bring

a holistic perspective that current AI techniques may lack.

An exciting development in the field of SLR is the relatively recent introduction of AI tools for automating the entire procedure — a field anticipated to continue expanding in the coming years. The increasing curiosity level indicates that now is an opportune time to analyze AI techniques presented as solutions to various SLR tasks. This analysis includes a focus on their intended use, sources of input and output, and the need for human intervention. Several research efforts in the field have incorporated AI techniques into their evaluations of procedures and instruments for facilitating SLR tasks. However, these investigations have taken either a more general approach, considering any type of automation with or without AI, or they have exclusively focused on AI [47], [76]. Some experts have concentrated on the use of specific AI techniques, such as ML methods, to address a particular problem [49] or a specific SLR activity, like document selection [57], [58].

Despite these efforts, some investigations may lack a comprehensive overview of the diverse ideas and procedures involved in AI relevant to the entire SLR procedure. In addition to providing a comprehensive overview of the field, the present paper aims to expand on the significance of human involvement—a perspective not fully addressed by the partially autonomous SLR considered in the existing literature reviews. Keeping these goals in mind, the following are a few inquiry concerns, also known as Research Questions (RQs), that inform our analysis of the current status of AI-based SLR automation:

- RQ1: Which stages of the SLR process have been automated using artificial intelligence?

- RQ2: Which AI methods facilitate the automation of SLR tasks?

- RQ3: To what extent does the human factor into SLR automation with AI?

As part of our survey, we conducted a systematic literature search to address the above RQs. We identified the latest original research articles from an extensive collection of references retrieved through both mechanical and human search systems. Reviewing these articles was essential to comprehend the motivation behind employing AI for specific tasks. We then scrutinized the inputs, outputs, and algorithmic choices of the proposed methods, along with information on the experimental evaluation of the approaches, including benchmarking metrics and sample articles.

Our analysis revealed that certain SLR tasks have been the subject of significantly more research compared to others, with some ML approaches introduced in the early phases still in use. However, we also identified more recent studies investigating novel ML approaches that incorporate the human dimension. Our findings in response to each RQ allowed us to pinpoint several unresolved concerns and difficulties related to the use of AI techniques for SLR tasks that they were not specifically designed for. Additionally, we identified issues related to experimental repeatability and other factors that have not yet been thoroughly addressed.

The remainder of this paper is organized as follows: Section II provides an overview of related work, highlighting existing literature and studies pertinent to the integration of AI in systematic literature reviews. In Section III, we delve into the methodology employed in our research, elucidating the approach and techniques used. Section IV explores the landscape of AI-based support for the literature review process, detailing advancements, tools, and strategies. Following this, Section V outlines open issues and challenges associated with AI-driven literature reviews. Section VI offers conclusions drawn from our exploration and proposes avenues for future research. These sections collectively contribute to a comprehensive understanding of the current state and potential future developments in the intersection of AI and systematic literature reviews.

## II. RELATED WORK

A systematic examination of existing research, known as a Systematic Literature Review (SLR), is a type of secondary investigation in a research field that systematically combines and evaluates scientific research to synthesize recent information, critically discuss current initiatives, and detect research patterns. SLRs use established procedures for conducting empirical research [33]. In particular, within Software Engineering (SE), researchers have made efforts to provide a comprehensive summary of methods devised for automating the SLR procedure. With a methodical approach to searching, they have reviewed the literature to shed light on different approaches used to automate various aspects of the SLR process [20].

In this context, our focus shifts to the work presented in [19], which demonstrates how computer languages can facilitate unsupervised ML for the synthesis and abstraction of data sets taken from an SLR. This article skillfully showcases the complementary roles that AI and ML techniques play in coding, categorization, and synthesis of SLR data, utilizing the qualitative method Deductive Qualitative Analysis [5].

While SLRs offer a clear and concise format for summarizing expertise in a field, they are not without challenges, such as the time required to complete them and the challenging task of assessing the integrity of primary research [35]. Recent analysis has highlighted prevalent hazards associated with SLR replication, emphasizing issues resulting from the absence of a defined methodology [38]. The approach [33] divides the SLR procedure into the following stages:

*1) Formulating phase:* The first aspect involves making a strategy. Justification for conducting an SLR in a research area ensures it addresses a gap and contributes to knowledge. Research queries are formulated to define the purview of the SLR and guide its evolution. These queries may adhere to predetermined structures, such as PICO (Population, Intervention, Comparison, and Outcome) or SPICE (Context, Perspective, Intervention, Comparison, and Evaluation) [15]. In this stage, an evaluation procedure is designed, including a comprehensive review technique applicable to each stage. The search technique and its sources, such as science resources and journals, are detailed in the protocol. Eligibility requirements for article selection, data extraction, and quality evaluation guidelines are also established.

*2) Conducting phase:* The second phase involves the execution of autonomous searches in data and digital libraries.

Search strings are obtained from either the formulated research queries or constructed using a supplementary method [50]. Additional sources, including dark texts and snowballs, are considered [42]. The former includes materials not publicly available, such as dissertations and presentations. The snowball effect involves discovering new literary works by examining references and citations from previously discovered papers. Relevant studies are identified by removing duplicates, evaluating candidates based on the name and summary, and applying inclusion and exclusion criteria. These criteria specify the quality standards each article must meet to be included in the scope [60]. The fundamental subjects are then analyzed to extract data, and summary statistics are obtained to synthesize and visualize the collected data.

*3) Reporting phase:* The third phase focuses on the reporting process and the evaluation of the final report's completeness and quality. Authors determine the manner in which the material is discussed and presented, as well as whether the evaluation result is suitable for publication. Criteria are considered to evaluate whether necessary data can be found in the SLR report [44].

## III. METHODOLOGY

### A. Search Strategy

The search strategy employs a combination of automated and human searching methods. Automated searches were conducted using the following sources: ACM Library, IEEE Xplore, Scopus, SpringerLink, and Web of Science. The search criteria, designed to retrieve publications, include a range of terms incorporating systematic review keywords and automation-related terms. General terms associated with automation were used rather than an exhaustive list of specific AI methods for two distinct purposes: (1) to avoid skewing the findings in favor of certain methods, ensuring inclusion of less prevalent ways in the final tally; and (2) to prevent the creation of lengthy and complicated search strings that may be challenging for databases to process. Title, keywords, and abstracts were considered in the search criteria. The resulting search string was easily adaptable for each data source. Additionally, a manual search was conducted using reverse snowballing. After reviewing the titles and abstracts of the initial eight candidate papers, six were added to the final list.

Fifty prospective papers were identified and underwent further evaluation to ensure their alignment with our research aims. For this purpose, both exclusion and inclusion criteria were developed. Papers written in languages other than English, those with unavailable full content, and publications lacking a demonstrated peer review process were excluded. The inclusion criteria set specific requirements for paper content. Each research paper must focus on automating multiple steps of an SLR and discuss the usage of AI-based methods for inclusion in the current survey. This general criterion is further subdivided into mutually exclusive options:

1) The paper explains a novel algorithm, instrument, or method facilitating full or partial mechanization of SLR;

2) It provides an examination of the relevance of AI in SLR, along with a critique of the latest developments in this field of study; and

3) The paper presents a summary of SLR tools applicable to one or more phases.

### B. Data Extraction

Once each primary research article has been identified, data extraction is performed following the guidelines outlined in [33]. One author reviews each article, with the assistance of a second reviewer in cases of ambiguity. The data extraction form includes meta-information such as authors and affiliations, research types, publication years, and publishing years.

The data extraction form also includes categories to define the AI approach followed in each paper. Specifically, each paper's content is summarized based on the following criteria:

*1) Phase and aim of the SLR:* Each paper is classified according to the phase of SLR automation, and each phase's categorization is followed by a description of the particular step(s) involved in that phase.

*2) AI domain and technique:* The paper is assigned to one or multiple automated academic subfields and contains a concise explanation of the employed algorithm or technique. We also record whether the societal factor is at play.

*3) Experimental framework:* Types of primary research include empirical, theoretical, application, and review. We compile the body of data and the indicators used for performance evaluation for empirical investigations.

*4) Repeatability:* Changes are made if any of the supplied tools, datasets, or algorithms require them. We verify the accessibility of any websites or repositories cited as supplementary material to ensure repeatability.

## IV. AI-BASED SUPPORT FOR THE LITERATURE REVIEW PROCESS

A literature review process often involves some visionary and mechanical constraints, prompting the development of AI-based technologies to alleviate the workload for potential authors. AI technologies aim to handle time-consuming and repetitive tasks, allowing authors to focus on interpretation, intuitive leaps, and skills [72].

To provide readers with insight into the current state of knowledge in this domain, we systematically evaluate each stage of the literature review process, highlighting existing AI-based tools and discussing the potential for further AI assistance. The following agenda can outline opportunities for additional industrial development and enhancement [69].

Table I presents a concise overview, indicating whether each stage is capable of being assisted by AI and guiding the reader toward relevant tools. The term "melted brief" refers to a succinct summary, capturing the essence of each stage's AI-assistive potential.

In particular, we conducted a comprehensive review of relevant literature on AI-based tools, consulting sources such as [1], [21], [26], [37]. Given the dynamic nature of AI technologies and the rapid evolution of AI tools, we focused

TABLE I. AI-BASED TOOLS FOR STEPS IN THE REVIEW PROCESS

| Steps | AI-based Tools | Potential for AI-support |
|---|---|---|
| 1. Problem Formulation | 1) Resources for software development supporting thematic analyses based on LDA models [3].<br>2) GUI applications and programming libraries supporting scientometric analyses [67]. | 1) Moderate potential with AI potentially pointing researchers to promising areas and questions or verifying research gaps. |
| 2. Literature Search | 1) TheoryOn enables ontology-based searches for constructs and construct relationships in behavioral theories [43].<br>2) Litbaskets supports researchers in setting a manageable breadth in regard to the magazines that have been covered [9].<br>3) LitSonar allows syntax interpretation of queries between database servers, as well as (journal coverage) creating reports [66]. | 1) Very high potential since the most important search methods consist of steps that are repetitive and time-consuming, that is, amenable to automation. |
| 3. Screening for Inclusion | 1) ASReview offers screening prioritization [75].<br>2) ADIT approach for researchers capable of designing and programming ML classifiers [40]. | 1) A significant opportunity for partly automated assistance in the initial stage screen, which requires many repetitive decisions.<br>2) Substantial possibility of an extra screen, requiring substantial expert judgement (particularly for ambiguous cases). |
| 4. Quality Assessment | 1) Statistical software packages (e.g., RevMan).<br>2) RobotReviewer for experimental research [48]. | 1) There is a possibility for partially automated quality control ranging from low to considerable. |
| 5. Data Extraction | 1) Software for data extraction and qualitative content analysis (e.g., Nvivo and ATLAS.ti) offers AI-based functionality for qualitative coding, named entity recognition, and sentiment analysis.<br>2) WebPlotDigitizer and Graph2Data for extracting data from statistical plots. | 1) Moderate potential for reviews requiring formal data extraction (descriptive reviews, scoping reviews, meta-analyses, and qualitative systematic reviews).<br>2) Elevated for quantitative and discrete data points (e.g., sample sizes), low for detailed information that is ambiguous and open to multiple interpretations (e.g., theorizing and main results). |
| 6. Data Analysis and Interpretation | 1) Descriptive synthesis Tools for text-mining [36], scientometric techniques, topic models [55], [63], and computational reviews aimed at stimulating conceptual contributions [4].<br>2) Theory building: Examples of inductive (computationally intensive) theory development [8], [45], [56].<br>3) Tools for doing meta-analyses, such as RevMan and dmetar, are used in the testing of hypotheses. | 1) Very high potential for descriptive syntheses.<br>2) Moderate potential for (inductive) theory development and theory testing.<br>3) Low non-existent potential for reviews adopting traditional and interpretive approaches. |

our evaluation on those most pertinent to our primary objective in the realm of Information Systems (IS) research. It is also important to note that our review is not exhaustive; rather, its purpose is to spotlight potential examples that can benefit IS researchers. In the upcoming paragraphs, we take a focused approach, examining AI-supported tasks individually. This approach allows us to provide insights into tools that authors can seamlessly integrate into a comprehensive data processing and toolchain.

### A. Step 1: Problem Formulation

In the initial phase of a comprehensive literature research, authors bear the responsibility of not only defining and elucidating research topics but also elucidating the fundamental ideas and theories within the relevant field [69]. Furthermore, scholars are advised to conduct a preliminary assessment of the research gap, determining whether the gap has been adequately addressed. Evaluating if the study's issue offers an opportunity for a substantial contribution that surpasses previous works and determining its significance in filling the existing void are crucial considerations in this phase [54], [62].

We envisage that AI can significantly contribute to the synthesis of research issues, particularly in the phase focused on identifying and validating open questions. With substantial advancements in the scientific domain, researchers have made significant strides in pinpointing gaps in the current body of knowledge and formulating plausible hypotheses. In this context, we anticipate that social science researchers can leverage and adapt these findings to enhance their own work. For instance, revolutionary developments in automated hypothesis generation and experimental testing have emerged in biochemistry, particularly within fully automated labs [34]. Additionally, ML strategies have been employed in scientometric approaches, facilitating literature-based discoveries in computer science [70]. These advancements underscore the potential for AI to play a transformative role in issue synthesis across various scientific disciplines.

Significant strides in information technology, particularly in the realm of database inquiries, have garnered attention. Three noteworthy resources underscore these recent advancements. Notably, TheoryOn's search engine [43] stands out. While these advancements hold promise, especially for studies in the

social sciences, their ultimate influence on research method-ologies remains to be seen. In general, these technological developments may prompt investigators to identify areas warranting further exploration or additional research. However, we acknowledge that, for problematization-driven research that generates queries, a continued reliance on deliberative democracy, particularly in the phase of problem identification, may be necessary [2].

In addition to recognizing voids within existing studies and areas calling for extensive exploration, AI holds the potential to assist researchers in assessing whether these gaps persist by locating prior assessments with similar or identical content. However, it's important to acknowledge that utilizing AI in this context may introduce a degree of unpredictability during the phases of discovery and verification, particularly when identifying prior assessments with content resembling the current study.

In conclusion, the support for this pioneering initiative in AI-backed research is still in its infancy, marked by a limited number of documented approaches. Notably, the existing software operates autonomously, diverging from the reliance on established graphical user interface tools. For researchers proficient in programming, inspiration can be drawn from exploring the intersections of various literature or study areas. This exploration opens avenues for cross-disciplinary research and identifies areas where further investigation is warranted. To facilitate this exploration, researchers can leverage scientometric approaches and enhance their development [18], [67]. Additionally, employing tools such as the LDA subject model can contribute to the identification of potential research directions [3]. As the field evolves, there is significant room for growth, and researchers are encouraged to embrace the interdisciplinary nature of AI-backed research to unlock its full potential.

*B. Step 2: Literature Search*

In this stage, researchers embark on constructing a comprehensive corpus of published works utilizing a diverse set of search techniques, including database queries, perusal of table of contents, citation inquiries, and additional searches [69]. The objective of the literature search can vary, with authors striving for either comprehensive, representative, or selective coverage based on the review's purpose [13].

Complex search strategies, involving multiple iterations and leveraging the collaboration with artificial intelligence-based technologies, are devised from corresponding search methods. Given the diverse nature of the knowledge retrieval process, encompassing various data sources such as journals, conference proceedings, books, and various forms of grey literature, alongside concerns about data quality, authors must employ appropriate data management strategies. These strategies should not only facilitate transparent reporting [61] but also contribute to repeatability and reproduction [14].

Recent advancements in information technology, particularly in the realm of database inquiries, have been noteworthy. Three notable resources stand out in this regard. Firstly, TheoryOn's search engine [43] enables researchers to conduct ontology-based inquiries for distinct elements and interactions between themes across different behavioral hypotheses. This

offers an alternative to traditional databases and presents a more sophisticated approach to information retrieval. Secondly, Litbaskets [9] contributes to the development of search methodologies by remotely estimating the likely number of responses from a database search based on predefined phrases across various journals with editable entries. Thirdly, LitSonar [66] adds automation to the search process by translating search terms for multiple book systems, including databases like EBSCO Digital Library, AIS eLibrary, and Lexisnexis. This tool is particularly promising as it provides real-time updates on service availability, potentially identifying database prohibitions (periods during which articles are not searchable) and alleviating challenges associated with database deficiencies.

In a broader context, the literature search phase holds the potential for automation, addressing numerous technological tasks that researchers encounter. The expanding volume of research output, coupled with the imperative need for efficiency and accuracy, underscores the significance of incorporating robotics and AI assistance in activities that are predominantly mechanical in nature [10], [25]. The integration of these technologies not only expedites the literature search process but also mitigates the inefficient use of faculty time resulting from manual tasks. This becomes especially critical in an era marked by the rapid proliferation of scholarly content.

*C. Step 3: Screening for Inclusion*

In this pivotal phase of the literature review, the authors employ a systematic screening process to differentiate between relevant and irrelevant papers. Conventionally, this phase is bifurcated into a preliminary assessment based on titles and abstracts, followed by a more rigorous second screening based on full-texts [69].

Manual screening, involving the meticulous examination of hundreds or thousands of documents, can be mentally taxing, potentially hindering the accurate identification of challenging scenarios. To mitigate this challenge, researchers are advised to conduct an initial screening where obviously irrelevant articles (based on titles and abstracts) are excluded. Articles posing difficulty are intentionally saved for a more comprehensive evaluation in the second round of screening.

The second screening involves a smaller sample, allowing for efficient screening (after excluding the majority in the initial screen). This stage involves a thorough examination of materials, application of predefined stringent exclusion criteria, and simultaneous independent evaluations, with group decisions on borderline cases. The screening process, particularly in hypothesis-testing reviews, requires stringent scrutiny, as inclusion errors could significantly impact the study outcomes [69].

Over time, the landscape of screening tools has evolved with the incorporation of AI-based tools [21]. Among them, ASReview [75], a tool with minimal limitations, stands out as a promising option for IS researchers. Operating on health sciences databases and requiring PubMed IDs are not significant limitations for ASReview. This tool, developed recently, is noteworthy for its transparency (under the Apache-2.0 License), script accessibility (implemented in Python), and the

ability to easily integrate new features. ASReview employs various ML classifiers, including Naive Bayes, logistic regression, Logical Regressive Analysis, and randomly generated forest classifiers. It leverages initial diversity decisions to enhance subsequent decision accuracy. Researchers receive a ranked catalog of papers (titles and descriptions), facilitating efficient processing of an ordered compilation. The tool even allows for automated exclusion after screening a specific number of papers consecutively, streamlining the screening process. Papers with borderline relevance can be deferred for later assessment, guided by their content [75].

For researchers with coding proficiency, customization of the discourse method is possible [40]. In situations where the evaluation of popular theories becomes impractical due to the sheer volume of pertinent documents, the prospect of randomly selecting theory-contributing publications is proposed. This approach leverages algorithms used in ML to identify a subset of relevant journals from the larger pool, thereby offering a randomized sample typical of how scientists articulate their work during literature reviews [40].

Anticipating the AI-support potential, the first screening demonstrates high efficacy, while the potential for the second screening is moderate. The initial screen, involving fewer exclusions, is more amenable to digitization and AI assistance. This presupposes computers with proficient reading and comprehension capabilities for brief descriptions and titles. Conversely, the subsequent screen deals with the remaining instances and may prove challenging due to the less standardized nature of IS research. Unlike fields like Medical and Biological Sciences, IS research lacks commonly used categories for constructs, standard keyword vocabulary (e.g., MeSH terms), and consistently descriptive paper titles, making effective classification challenging [58]. This challenge is not exclusive to machines and equally affects human reviewers.

Screening and search, treated as information retrieval tasks, should primarily be evaluated based on recall, representing the proportion of successfully retrieved relevant papers. Traditionally, literature reviews aimed for high recall, resulting in exhaustive searches, low precision, and increased screening burdens [43]. AI-supported ontology-based searches, such as those facilitated by ASReview, hold the promise of efficiently alleviating a portion of the screening load by increasing precision.

The screening processes remain among the most time-consuming aspects of a literature review process [10]. When considering the potential of AI assistance for these steps, it's crucial to recognize that the reliability of manual screening methods should not be overstated, as even screenings conducted by experts exhibit a disagreement rate of 10% on average [81]. Augmenting researchers' screening activities with AI tools can help identify inconsistent and potentially erroneous screening decisions, enhancing the reliability of the screening process.

### D. Step 4: Quality Assessment

The evaluation of major empirical research for methodological flaws and potential sources of bias is an integral part of the quality assessment process [23], [33], [69]. This phase aims to gauge the extent to which the findings of

evaluations, particularly those intended for theory testing, may be influenced by various types of bias, such as selection bias, mortality bias, and evaluation bias. Parallel and independent execution of these procedures is recommended to ensure high dependability [69].

The prospect of AI-based tools contributing to these processes is considered to have a low to middling chance for two primary reasons. Firstly, the task of judging the quality of a method is challenging, requiring expert opinion and often presenting difficulties in achieving high inter-coder agreement [22]. Secondly, IS reviews, whether quantitative or qualitative, typically involve manageable numbers of samples, making manual assessments feasible.

For researchers conducting meta-analyses and systematic literature searches, conventional tools like RevMan, adhering to standards for evaluating qualitative research methodology and risk of bias [7], or equivalent statistical application environments such as R and SPSS are commonly used. Additionally, AI-based applications like RobotReviewer [48] offer relevance to IS meta-analyses. RobotReviewer, focusing on risk of bias assessment in randomized controlled trials within the life sciences, serves as an exemplary instance of explainable AI. It enables scholars to trace ratings in each bias area back to their source within the full-text document. This transparency contributes to the reliability and interpretability of the bias assessment process.

### E. Step 5: Data Extraction

The extraction of data, both qualitative and quantitative, involves the identification of relevant information and its categorization into a (semi) structured code sheet [69], [79]. This step is more prominent in description, scoping, and theory testing reports compared to narrative reviews and assessments of theoretical development, which tend to be more selective and interpretive. Commonly utilized software for all-encompassing subjective data analysis includes ATLAS.ti and NVivo, which are increasingly incorporating ML and NLP techniques. These techniques include methods for information extraction from data tables, automation of descriptive coding, Named Entity Recognition (NER), sentiment analysis, and analysis of statistical plots. Examples of tools for extracting data from statistical plots include WebPlotDigitizer and Graph2Data.

The potential for AI support in this step is anticipated to be moderate. Future advancements may focus on improving the efficiency of information extraction, highlighting crucial elements in an article, and facilitating the organization of information in suitable databases. However, complete automation of more intricate data elements is not expected in the near future. Despite the more standardized disclosure practices in the medical professions, tools for extracting features such as Population, Intervention, Comparison, and Outcome (PICO) criteria are still in their early stages of development [26].

### F. Step 6: Data Analysis and Interpretation

The concluding stage of the evaluation process in literature reviews can take various forms depending on the type of assessment [69]. Some literature reviews emphasize intricate scenarios that provide insights and profound hermeneutic interpretations, while others aim to eliminate subjectivity that

might compromise the reliability of summary statistics and generalizations.

IS researchers employ various instruments for data analysis, depending on the main objectives for knowledge development [64]. For comprehensive synthesis, several well-established techniques are available, including text-mining tools [36] and instruments that utilize scientometric, computational, or Latent Dirichlet Allocation (LDA) models to analyze and visualize themes, theories, and research communities [6], [17], [41], [55], [68], [70], [77], [78]. For example, text-mining tools can provide descriptive insights based on topic modeling, offering a promising approach to conceptual contributions [39], [53], [63]. In the realm of IS, meta-analysis programs and libraries, such as RevMan and the R package dmetar [7], are utilized for putting hypotheses to the test. Future AI-based technologies supporting data analysis should consider the diverse approaches available. While AI can efficiently ease certain aspects of descriptive evaluations through topic modeling, the creative and unstructured nature of theory development poses challenges for AI-led theory-building efforts [8], [45], [56].

The inductive method of IS theory development seems most amenable to AI assistance, although current examples in the behavioral research domain may not match the ingenuity and originality exhibited by exceptional theoretical and historical context articles [4], [39], [53], [63]. It is crucial to emphasize that AI's contribution to theory development lies not only in identifying connections but also in elucidating the "why" behind these connections and establishing fundamental philosophical foundations of justification [25], [82]. This aspect remains an unrestricted challenge for future theory advancement based on AI.

Fig. 1 illustrates the three layers of the SLR-centric approach, emphasizing the goals for research, design, and action within the Information Systems (IS) domain. This three-layered SLR-centric model within the IS domain underscores the interplay between infrastructure, methodologies/tools, and actual research practices, showcasing the integral components for successful literature reviews.



Fig. 1. SLR-Centric research, design, and action.

Quality assurance is integral to ensuring the accuracy, dependability, and consistency of data collected and analyzed

during the literature review [23], [33], [69]. Smart technologies, incorporating automation and clever algorithms, enhance the efficiency and effectiveness of literature review operations [30], [31], [46], [59], [73], [80]. Improved databases are essential for efficiently storing, maintaining, and retrieving literature review data [14]. Integrating these features into the infrastructure supporting SLR substantially enhances the quality, efficiency, and insights of the review process.

The standardization debate in the realm of Information Systems (IS) revolves around whether standardized approaches, methodologies, and technology should be embraced across multiple IS applications or whether variety and flexibility should be prioritized [32]. The concept of sharing supplementary research outputs emphasizes the necessity of sharing not only traditional research publications but also other significant outputs contributing to the research process. This includes datasets, code, methodology, negative outcomes, and other relevant items [64].

## V. Open Issues and Challenges

### A. The Emphasis is Primarily on One Activity

Automating SLR processes using AI research is heavily skewed towards the execution of paper selection procedure, especially during the phase of paper assignment. While this activity is also time-consuming, it is important to note that applying AI to various tasks within the SLR process requires attention. Preliminary tasks, such as AI-driven writing activities (e.g., drafting research questions, specifying exclusion/inclusion criteria, and presenting SLR reports), are areas that need further development.

### B. More Research Needs to Be Done on AI Methods

While there is a broad range of AI fields and techniques, certain ones have not yet been employed in SLR automation. Methods for optimization and inquiry, for instance, have not been thoroughly investigated as potential solutions for SLR-related tasks. These strategies, traditionally used to resolve planning issues, may find application in prioritizing resources during the initial stages, such as selecting the best databases or assigning papers to editors based on their skills. Compared to ML, knowledge representation and NLP are less frequent, and most proposals appear to be in early stages. Therefore, there is a need for more tools and frameworks to develop solutions based on these methods.

### C. Additional Active Human Participation can Benefit Artificial Intelligence

The cooperation between humans and AI methods or instruments is currently limited in terms of scope and nature. Under an active learning strategy, the human role is primarily focused on providing labels for paper selection. However, the organizing and composing phases, which demand greater human capabilities, might benefit from engaging artificial intelligence. Involving people in this process could result in additional positive outcomes, such as tailoring the results to their preferences.

## D. The Adoption of AI for SLR Automation can be Enhanced

Most current successful ideas are rooted in either the medical or technological domains, with specific domain taxonomies or concepts sometimes used to construct the list of capabilities. Full replicability of genuine systematic literature reviews is not always achieved, and the lack of benchmarks remains a significant obstacle. Evaluating AI techniques across a broader range of SLRs and expanding the scope of discussed issues requires additional development.

## E. Users of SLR Automation may Lack Expertise in Artificial Intelligence

Many ML techniques examined thus far, such as support vector machines (SVM) and neural networks, which are commonly referred to as "black-box" techniques. The challenge arises from insufficient confidence in automated conclusions due to the participation of scientists from diverse disciplines in SLRs, who may not necessarily be experts in AI. There has been limited exploration of models using human-readable code, including simple decision trees and rule-based systems. Furthermore, the utilization of contemporary explainable techniques has the potential to enhance the outcomes of black-box artificial intelligence solutions developed in this field.

## VI. CONCLUSIONS AND FUTURE WORK

### A. Conclusions

Literature reviews represent just one facet of a laborious and error-prone process that can be streamlined with the assistance of artificial intelligence (AI). It is not surprising that not all tasks involved in Systematic Literature Review (SLR) planning, execution, and reporting have been fully automated to date. Our research indicates a strong inclination to leverage AI, particularly Machine Learning (ML), to facilitate the paper screening process, which entails sifting through thousands of candidate papers to identify relevant ones. Natural Language Processing (NLP) and ontologies prove especially beneficial in handling semantic data for various tasks, although there is a paucity of studies in these domains.

The findings highlight the need for a strategic roadmap for AI-led research, development, and implementation across different dimensions. Our primary objective is to foster the growth of a vibrant AI-based Literature Reviews (AILR) culture in the Information Systems (IS) field, offering enriching experiences for researchers at all stages of the research process—from authors to reviewers to industry professionals. The potential for extending AI-based tools and approaches beyond the IS field, particularly for design science researchers, holds great promise. We envision a future where IS researchers actively engage in discussions and reflections on how to optimally harness AI to advance their work.

### B. Future Work

The future work outlined in the provided passage encompasses a comprehensive and forward-thinking strategy for advancing the field of AI-based literature reviews. Initially, the focus is directed towards extending the application of AI across various phases of the systematic literature review (SLR) process. While the current emphasis is on paper screening, the call is for a more encompassing integration, suggesting a desire to leverage AI capabilities throughout the entire SLR workflow. This expansion could lead to more nuanced and sophisticated automation, enhancing the efficiency and effectiveness of literature review processes.

Additionally, the passage underscores the importance of exploring advanced technologies, such as Natural Language Processing (NLP) and ontologies, for semantic data analysis in the context of literature reviews [16]. This suggests an aspiration to move beyond basic automation and delve into the realms of semantic understanding, potentially enabling AI systems to discern and interpret the underlying meaning of academic content. Moreover, the envisioned future involves broadening the application of AI-based tools and methodologies to domains beyond Information Systems (IS), emphasizing the need for transferability and interdisciplinary collaboration. This push for broader applicability aligns with the broader trend of fostering cross-disciplinary knowledge exchange and collaboration.

Furthermore, the future work envisions a significant enhancement of the user experience for researchers engaging with AI tools in the SLR process. This user-centric approach aims to make AI more accessible to individuals with varying skill levels, fostering inclusivity and democratizing the use of advanced technologies in academia. Simultaneously, there is a recognition of the importance of establishing benchmarks and evaluation criteria to assess the effectiveness and efficiency of AI-driven SLR processes. This emphasis on standardization reflects a commitment to ensuring robust and comparable outcomes in the application of AI methodologies to literature reviews. Finally, ethical considerations emerge as a crucial aspect of the envisioned future work, with a call to address ethical implications and develop guidelines for responsible AI implementation in research processes. This reflects a conscientious approach towards deploying AI in a manner that upholds ethical standards and promotes responsible conduct in academic research. In summary, the future work outlined in the passage is characterized by a holistic vision, encompassing technological advancements, usability improvements, interdisciplinary applications, ethical considerations, and a commitment to standardization in the realm of AI-based literature reviews.

## REFERENCES

[1] A. Al-Zubidy, J. C. Carver, D. P. Hale, and E. E. Hassler. Vision for slr tooling infrastructure: Prioritizing value-added requirements. *Information and Software Technology*, 91:72–81, 2017.

[2] M. Alvesson and J. Sandberg. Generating research questions through problematization. *Academy of management review*, 36(2):247–271, 2011.

[3] D. Antons and C. F. Breidbach. Big data, big insights? Advancing service innovation and design with machine learning. *Journal of Service Research*, 21(1):17–39, 2018.

[4] D. Antons, C. F. Breidbach, A. M. Joshi, and T. O. Salge. Computational literature reviews: Method, algorithms, and roadmap. *Organizational Research Methods*, 26(1):107–138, 2023.

[5] C. F. Atkinson. Cheap, quick, and rigorous: Artificial intelligence and the systematic literature review. *Social Science Computer Review*, page 08944393231196281, 2023.

[6] B. Balducci and D. Marinova. Unstructured data in marketing. *Journal of the Academy of Marketing Science*, 46:557–590, 2018.

[7] L. Bax, L.-M. Yu, N. Ikeda, and K. G. Moons. A systematic comparison of software dedicated to meta-analysis of causal studies. *BMC medical research methodology*, 7:1–9, 2007.

[8] N. Berente, S. Seidel, and H. Safadi. Research commentary—data-driven computationally intensive theory development. *Information Systems Research*, 30(1):50–64, 2019.

[9] S. Boell and B. Wang. An it artifact supporting exploratory literature searches. In *Australasian conference on information systems. http://www.litbaskets.io. Accessed*, volume 21, 2021.

[10] J. C. Carver, E. Hassler, E. Hernandes, and N. A. Kraft. Identifying barriers to the systematic literature review process. In *2013 ACM/IEEE international symposium on empirical software engineering and measurement*, pages 203–212. IEEE, 2013.

[11] A. L. Chapman, L. C. Morgan, and G. Gartlehner. Semi-automating the manual literature search for systematic reviews increases efficiency. *Health Information & Libraries Journal*, 27(1):22–27, 2010.

[12] A. M. Cohen, W. R. Hersh, K. Peterson, and P.-Y. Yen. Reducing workload in systematic review preparation using automated citation classification. *Journal of the American Medical Informatics Association*, 13(2):206–219, 2006.

[13] H. M. Cooper. Organizing knowledge syntheses: A taxonomy of literature reviews. *Knowledge in society*, 1(1):104, 1988.

[14] W. A. Cram, M. Templier, and G. Paré. (re) considering the concept of literature review reproducibility. *Journal of the Association for Information Systems*, 21(5):10, 2020.

[15] K. S. Davies. Formulating the evidence based practice question: a review of the frameworks. *Evidence Based Library and Information Practice*, 6(2):75–80, 2011.

[16] G. Drakopoulos, A. Kanavos, P. Mylonas, S. Sioutas, and D. Tsolis. Towards a framework for tensor ontologies over neo4j: Representations and operations. In *8th International Conference on Information, Intelligence, Systems & Applications (IISA)*, pages 1–6. IEEE, 2017.

[17] E. Dritsas, M. Trigka, G. Vonitsanos, A. Kanavos, and P. Mylonas. Aspect-based community detection of cultural heritage streaming data. In *12th International Conference on Information, Intelligence, Systems & Applications (IISA)*, pages 1–4. IEEE, 2021.

[18] J. A. Evans and J. G. Foster. Metaknowledge. *Science*, 331(6018):721–725, 2011.

[19] K. R. Felizardo and J. C. Carver. Automating systematic literature review. *Contemporary empirical methods in software engineering*, pages 327–355, 2020.

[20] K. R. Felizardo, É. F. de Souza, B. M. Napoleão, N. L. Vijaykumar, and M. T. Baldassarre. Secondary studies in the academic context: A systematic mapping and survey. *Journal of Systems and Software*, 170:110734, 2020.

[21] H. Harrison, S. J. Griffin, I. Kuhn, and J. A. Usher-Smith. Software tools to support title and abstract screening for systematic reviews in healthcare: an evaluation. *BMC medical research methodology*, 20:1–12, 2020.

[22] L. Hartling, M. Ospina, Y. Liang, D. M. Dryden, N. Hooton, J. K. Seida, and T. P. Klassen. Risk of bias versus quality assessment of randomised controlled trials: cross sectional study. *Bmj*, 339, 2009.

[23] J. P. Higgins, J. Thomas, J. Chandler, M. Cumpston, T. Li, M. J. Page, and V. A. Welch. *Cochrane handbook for systematic reviews of interventions*. John Wiley & Sons, 2019.

[24] M.-S. James, C. Marrissa, S. Mark, B. Anthea, et al. Systematic approaches to a successful literature review. *Systematic Approaches to a Successful Literature Review*, pages 1–100, 2021.

[25] C. D. Johnson, B. C. Bauer, and F. Niederman. The automation of management and business science. *Academy of Management Perspectives*, 35(2):292–309, 2021.

[26] S. R. Jonnalagadda, P. Goyal, and M. D. Huffman. Automating data extraction in systematic reviews: a systematic review. *Systematic reviews*, 4(1):1–16, 2015.

[27] A. Kanavos, N. Antonopoulos, I. Karamitsos, and P. Mylonas. A comparative analysis of tweet analysis algorithms using natural language processing and machine learning models. In *18th International Workshop on Semantic and Social Media Adaptation and Personalization (SMAP)*, pages 1–6. IEEE, 2023.

[28] A. Kanavos, C. Makris, Y. Plegas, and E. Theodoridis. Ranking web search results exploiting wikipedia. *International Journal on Artificial Intelligence Tools*, 25(3):1650018:1–1650018:26, 2016.

[29] A. Kanavos, E. Theodoridis, and A. K. Tsakalidis. Extracting knowledge from web search engine results. In *24th International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 860–867. IEEE Computer Society, 2012.

[30] I. Karamitsos, M. Papadaki, and N. B. Al Barghuthi. Design of the blockchain smart contract: A use case for real estate. *Journal of Information Security*, 9(3):177–190, 2018.

[31] I. Karamitsos, M. Papadaki, K. Al-Hussaeni, and A. Kanavos. Transforming airport security: Enhancing efficiency through blockchain smart contracts. *Electronics*, 12(21):4492, 2023.

[32] I. Karydis, A. Kanavos, S. Sioutas, M. Avlonitis, and N. I. Karacapilidis. Multimedia content's brokerage: An information system based on lesim. *International Journal of E-Services and Mobile Applications*, 12(2):40–58.

[33] S. Keele et al. Guidelines for performing systematic literature reviews in software engineering, 2007.

[34] R. D. King, J. Rowland, S. G. Oliver, M. Young, W. Aubrey, E. Byrne, M. Liakata, M. Markham, P. Pir, L. N. Soldatova, et al. The automation of science. *Science*, 324(5923):85–89, 2009.

[35] B. Kitchenham and P. Brereton. A systematic review of systematic review process research in software engineering. *Information and software technology*, 55(12):2049–2075, 2013.

[36] V. B. Kobayashi, S. T. Mol, H. A. Berkers, G. Kismihók, and D. N. Den Hartog. Text mining in organizational research. *Organizational research methods*, 21(3):733–765, 2018.

[37] C. Kohl, E. J. McIntosh, S. Unger, N. R. Haddaway, S. Kecke, J. Schiemann, and R. Wilhelm. Online tools supporting the conduct and reporting of systematic reviews and systematic maps: a case study on cadima and review of existing tools. *Environmental Evidence*, 7:1–17, 2018.

[38] J. Krüger, C. Lausberger, I. von Nostitz-Wallwitz, G. Saake, and T. Leich. Search. review. repeat? An empirical study of threats to replicating slr searches. *Empirical Software Engineering*, 25:627–677, 2020.

[39] M. Kunc, M. J. Mortenson, and R. Vidgen. A computational literature review of the field of system dynamics from 1974 to 2017. *Journal of Simulation*, 12(2):115–127, 2018.

[40] K. R. Larsen, D. Hovorka, A. Dennis, and J. D. West. Understanding the elephant: The discourse approach to boundary identification and corpus construction for theory review articles. *Journal of the association for information systems*, 20(7):15, 2019.

[41] C. Laurell, C. Sandström, A. Berthold, and D. Larsson. Exploring barriers to adoption of virtual reality through social media analytics and machine learning–an assessment of technology, network, price and trialability. *Journal of Business Research*, 100:469–474, 2019.

[42] C. Lefebvre, E. Manheimer, and J. Glanville. Searching for studies. *Cochrane handbook for systematic reviews of interventions: Cochrane book series*, pages 95–150, 2008.

[43] J. Li, K. Larsen, and A. Abbasi. Theoryon: A design framework and system for unlocking behavioral knowledge through ontology learning. *MIS Quarterly*, 44(4), 2020.

[44] A. Liberati, D. G. Altman, J. Tetzlaff, C. Mulrow, P. C. Gøtzsche, J. P. Ioannidis, M. Clarke, P. J. Devereaux, J. Kleijnen, and D. Moher. The prisma statement for reporting systematic reviews and meta-analyses of studies that evaluate health care interventions: explanation and elaboration. *Annals of internal medicine*, 151(4):W–65, 2009.

[45] A. Lindberg. Developing theory through integrating human and machine pattern recognition. *Journal of the Association for Information Systems*, 21(1):7, 2020.

[46] A. K. Lingaraju, M. Niranjanamurthy, P. Bose, B. Acharya, V. C. Gerogiannis, A. Kanavos, and S. Manika. Iot-based waste segregation with location tracking and air quality monitoring for smart cities. *Smart Cities*, 6(3):1507–1522, 2023.

[47] C. Marshall, P. Brereton, and B. Kitchenham. Tools to support systematic reviews in software engineering: a feature analysis. In *Proceedings of the 18th international conference on evaluation and assessment in software engineering*, pages 1–10, 2014.

[48] I. J. Marshall, J. Kuiper, and B. C. Wallace. Robotreviewer: evaluation of a system for automatically assessing bias in clinical trials. *Journal of the American Medical Informatics Association*, 23(1):193–201, 2016.

[49] I. J. Marshall and B. C. Wallace. Toward systematic review automation: a practical guide to using machine learning tools in research synthesis. *Systematic reviews*, 8:1–10, 2019.

[50] G. D. Mergel, M. S. Silveira, and T. S. da Silva. A method to support search string building in systematic literature reviews through visual text mining. In *Proceedings of the 30th annual ACM symposium on applied computing*, pages 1594–1601, 2015.

[51] A. Mohasseb, B. Aziz, and A. Kanavos. SMS spam identification and risk assessment evaluations. In *16th International Conference on Web Information Systems and Technologies (WEBIST)*, pages 417–424, 2020.

[52] A. Mohasseb, M. Bader-El-Den, A. Kanavos, and M. Cocea. Web queries classification based on the syntactical patterns of search types. In *19th International Conference on Speech and Computer (SPECOM)*, volume 10458 of *Lecture Notes in Computer Science*, pages 809–819. Springer, 2017.

[53] M. J. Mortenson and R. Vidgen. A computational literature review of the technology acceptance model. *International Journal of Information Management*, 36(6):1248–1259, 2016.

[54] C. Müller-Bloch and J. Kranz. A framework for rigorously identifying research gaps in qualitative literature reviews. 2015.

[55] S. Nakagawa, G. Samarasinghe, N. R. Haddaway, M. J. Westgate, R. E. O'Dea, D. W. Noble, and M. Lagisz. Research weaving: visualizing the future of research synthesis. *Trends in ecology & evolution*, 34(3):224–238, 2019.

[56] L. K. Nelson. Computational grounded theory: A methodological framework. *Sociological Methods & Research*, 49(1):3–42, 2020.

[57] B. K. Olorisade, E. de Quincey, P. Brereton, and P. Andras. A critical analysis of studies that address the use of text mining for citation screening in systematic reviews. In *Proceedings of the 20th international conference on evaluation and assessment in software engineering*, pages 1–11, 2016.

[58] A. O'Mara-Eves, J. Thomas, J. McNaught, M. Miwa, and S. Ananiadou. Using text mining for study identification in systematic reviews: a systematic review of current approaches. *Systematic reviews*, 4(1):1–22, 2015.

[59] T. Panagiotakopoulos, D. P. Vlachos, T. V. Bakalakos, A. Kanavos, and A. Kameas. A fiware-based iot framework for smart water distribution management. In *12th International Conference on Information, Intelligence, Systems & Applications (IISA)*, pages 1–6. IEEE, 2021.

[60] D. Papaioannou, A. Sutton, and A. Booth. Systematic approaches to a successful literature review. *Systematic approaches to a successful literature review*, pages 1–336, 2016.

[61] G. Paré, M. Tate, D. Johnstone, and S. Kitsiou. Contextualizing the twin concepts of systematicity and transparency in information systems literature reviews. *European Journal of Information Systems*, 25:493–508, 2016.

[62] S. Rivard. Editor's comments: The ions of theory construction. 2014.

[63] T. Schmiedel, O. Müller, and J. Vom Brocke. Topic modeling as a strategy of inquiry in organizational research: A tutorial with an application example on organizational culture. *Organizational Research Methods*, 22(4):941–968, 2019.

[64] G. Schryen, G. Wagner, A. Benlian, and G. Paré. A knowledge development perspective on literature reviews: Validation of a new typology in the is field. *Communications of the AIS*, 46, 2020.

[65] C. Stansfield, A. O'Mara-Eves, and J. Thomas. Text mining for search term development in systematic reviewing: A discussion of some

methods and challenges. *Research synthesis methods*, 8(3):355–365, 2017.

[66] B. Sturm and A. Sunyaev. Design principles for systematic search systems: a holistic synthesis of a rigorous multi-cycle design science research journey. *Business & Information Systems Engineering*, 61:91–111, 2019.

[67] D. R. Swanson and N. R. Smalheiser. An interactive system for finding complementary literatures: a stimulus to scientific discovery. *Artificial intelligence*, 91(2):183–203, 1997.

[68] W. L. Tate, L. M. Ellram, and J. F. Kirchoff. Corporate social responsibility reports: a thematic analysis related to supply chain management. *Journal of supply chain management*, 46(1):19–44, 2010.

[69] M. Templier and G. Pare. Transparency in literature reviews: an assessment of reporting practices across review types and genres in top is journals. *European Journal of Information Systems*, 27(5):503–550, 2018.

[70] M. Thilakaratne, K. Falkner, and T. Atapattu. A systematic review on literature-based discovery: general overview, methodology, & statistical analysis. *ACM Computing Surveys (CSUR)*, 52(6):1–34, 2019.

[71] M. Torres Torres and C. E. Adams. Revmanhal: towards automatic text generation in systematic reviews. *Systematic reviews*, 6:1–7, 2017.

[72] G. Tsafnat, P. Glasziou, M. K. Choong, A. Dunn, F. Galgani, and E. Coiera. Systematic review automation technologies. *Systematic reviews*, 3:1–15, 2014.

[73] G. Tsaramirsis, I. Karamitsos, and C. Apostolopoulos. Smart parking: An iot application for smart city. In *2016 3rd International conference on computing for sustainable global development (INDIACom)*, pages 1412–1416. IEEE, 2016.

[74] A. Van Altena, R. Spijker, and S. Olabarriaga. Usage of automation tools in systematic reviews. *Research synthesis methods*, 10(1):72–82, 2019.

[75] R. Van De Schoot, J. De Bruin, R. Schram, P. Zahedi, J. De Boer, F. Weijdema, B. Kramer, M. Huijts, M. Hoogerwerf, G. Ferdinands, et al. An open source machine learning framework for efficient and transparent systematic reviews. *Nature machine intelligence*, 3(2):125–133, 2021.

[76] R. van Dinter, B. Tekinerdogan, and C. Catal. Automation of systematic literature reviews: A systematic literature review. *Information and Software Technology*, 136:106589, 2021.

[77] W. van Zoonen and G. Toni. Social media research: The application of supervised machine learning in organizational communication research. *Computers in human behavior*, 63:132–141, 2016.

[78] G. Vonitsanos, A. Kanavos, A. Mohasseb, and D. Tsolis. A nosql approach for aspect mining of cultural heritage streaming data. In *10th International Conference on Information, Intelligence, Systems and Applications (IISA)*, pages 1–4. IEEE, 2019.

[79] G. Vonitsanos, A. Kanavos, P. Mylonas, and S. Sioutas. A nosql database approach for modeling heterogeneous and semi-structured information. In *9th International Conference on Information, Intelligence, Systems and Applications (IISA)*, pages 1–8. IEEE, 2018.

[80] G. Vonitsanos, T. Panagiotakopoulos, A. Kanavos, and A. K. Tsakalidis. Forecasting air flight delays and enabling smart airport services in apache spark. In *Artificial Intelligence Applications and Innovations (AIAI)*, volume 628 of *IFIP Advances in Information and Communication Technology*, pages 407–417. Springer, 2021.

[81] Z. Wang, T. Nayfeh, J. Tetzlaff, P. O'Blenis, and M. H. Murad. Error rates of human reviewers during abstract screening in systematic reviews. *PloS one*, 15(1):e0227742, 2020.

[82] D. A. Whetten. What constitutes a theoretical contribution? *Academy of management review*, 14(4):490–495, 1989.

# Reciprocal Bucketization (RB) - An Efficient Data Anonymization Model for Smart Hospital Data Publishing

Rajesh S M[1], Prabha R[2]

Research Scholar, Department of Information Science and Engineering,
Dr. Ambedkar Institute of Technology, Visvesvaraya Technological University
Department of Computer Science and Engineering,
GITAM Deemed to be University, Bengaluru, India.[1]
Department of Computer Science and Engineering, Dr. Ambedkar Institute of Technology, Bengaluru, India.[2]

*Abstract*—With the lightning growth of the Internet of Things (IoT), enormous applications have been developed to serve industries, the environment, society, etc. Smart Health care is one of the significant applications of the IoT, where intelligent environments enrich safety and ease of surveillance. The database of the Smart Hospital records the patient's sensitive information, which could face various potential privacy breaches through linkage attacks. Publishing such sensitive data to society is challenging in adopting the best privacy preservation model to defend against linkage attacks. In his paper, we propose a novel Reciprocal Bucketization Anonymization model as the privacy preservation method to defend against Identity, Attribute, and Correlated Linkage attacks. The proposed anonymization method creates the Buckets of patient records and then partitions the data into sensor trajectory and Multiple Sensitive attributes (MSA). A local suppression is employed on Sensor Trajectory Data and Slicing on MSA to get the anonymized data to be published gathered by combining anonymized sensor trajectory and MSA. The proposed method is validated on the synthetic and real-time dataset by comparing its data utility loss in both sensor trajectory and the MSA. The experimental results eradicate that the RB – Anonymization exhibits the nature of best privacy preservation against Identity, Attribute, and Correlated linkages attacks with negligible utility loss compared with the existing methods.

*Keywords*—*Anatomization; anonymization; entropy; pearson's contingency coefficient; and KL – Divergence*

## I. INTRODUCTION

The Internet of Things (IoT) ecosystem facilitates collaboration across various computing devices ranging from sensors to complex processing systems with cloud storage. In the digitally connected world era, IoT adoption has rapidly increased in multiple applications, such as Smart Homes, Smart Healthcare, and so on [1]. The IoT plays a crucial role in building a secure cyber-physical communication system as the essential requirement for deploying Intelligent applications. Smart Healthcare systems primarily focus on patient real-time monitoring through sensors and data management via the cloud for remote access, such as Mobi-Care and MEDiSN [2].



Fig. 1. Structure of smart hospital.

Fig. 1 describes the components deployed in the Smart Hospital Structure with its various benefits to assess the patients and the caretakers to build a pervasive surveillance system. The evolution in the sensor's technology assures innovative and secure patient care in real-time and facilitates the accessible collection of the patient's spatiotemporal sensor data [3]. The sequence of sensor data at the specific time of a patient is known as sensor trajectory data. It helps to predict patient-sensitive information such as symptoms, diseases, etc. However, publishing such trajectory data along with multiple sensitive attributes of a patient for the researchers/data miners may result in a privacy breach [4]. Hence the challenge is preserving the privacy of the patient data by providing equal weightage for trajectory data and multiple sensitive attributes against the potential privacy breach.

Most privacy preservation principles are designed to protect the data against privacy breaches of patient sensitive information. However, the Adversary with prior knowledge of partial sensor trajectory data or a few sensitive attributes could infer the patient data even after removing the identity attributes, such as patient name and unique social identification number from the database [5]. Further, the Adversary can imply various linkage attacks with the prior knowledge to predict the patient's sensitive information with high probability [6].

Example: Consider the SMART Hospital "X," which digitally records and maintains patient data. The records may have the patient's ID, sensor trajectory, and medical data, as represented in Table I. The sensor trajectory data is recorded concerning time from the sensors deployed on the patient's

body, representing the pair of sensor data and time given as (sen,t) [7] [8]. An example of the Patient data is as follows, PId 7 is the patient ID, the data collected from the sensors x, k, n, and p at the timestamps 1, 6, 7, and 8, respectively, with the multiple sensitive attributes Fever, RITD Test, Influenza, and Medicine. The recorded data must be made available for the data miners for research [9]. In parallel, the hospital wants to preserve the privacy of the patient's sensitive data from unauthorized usage by malicious data miners against the following attacks [10]:

Identity linkage attack: In the published dataset, if the trajectory of the sensor data for a patient is unique, then an adversary can quickly identify a patient record along with the patient's sensitive data using his prior knowledge [11].

Attribute linkage attack: In the published dataset, the most frequent occurrence of the sensitive values of a targeted victim could result in an attribute linking attack. The adversary could breach the sensitive information with high confidence even though the unique sensor trajectory information of the victim is not available [11].

Correlated-records linkage attack: In the published dataset, when a patient has multiple records, there could be the possibility of a Correlated-record linkage attack. For example, the patient with PId 1 has the correlated records in row 1 and row 4, as shown in Table I. Having additional knowledge about the correlated records by the adversary can predict both trajectory data and the sensitive value of the victim with high confidence [11].

In literature, various privacy preservation approaches have been proposed for the trajectory data with single sensitive attributes, such as Generalization, Perturbation, Clustering, Differential Privacy, and Suppression, to defend against the various linkage attacks [12] [13]. Similarly, Multi-sensitive Bucketization, (p,k) —Angelization, and Generalization are the approaches to preserve the privacy of the multiple sensitive attributes along with the Quasi Identifiers but not with trajectory data. To the best of our knowledge, we are the first to address the privacy preservation approach for the dynamic trajectory data generated by the sensors and Multiple Sensitive attributes with trustfulness.

In this paper, we implement Reciprocal Bucketization as the overall framework for anonymization. Further Suppression and Slicing are implemented to anonymize the trajectory data with Multiple Sensitive Attributes to ensure privacy from the above three linkage attacks. The Bucketization approach helps in the formation of buckets from the patient data table records, on which the suppression and the slicing methods are parallelly imposed on sensor trajectory data and MSA, respectively, to anonymize the data via K-anonymity threshold by reducing the data loss with efficient anonymization [14] [15].

The major contributions are summarized as follows:

- We present Reciprocal Bucketization (RB) an efficient data anonymization model to preserve privacy in publishing the patient's data by the Smart Hospitals.

- To the best of our knowledge, we are the first to combine the Sensor trajectory data and MSA to ensure the privacy requirements for data publishing to defend against Identity, Attribute, and Correlated Linkage attacks.

- We proposed a suppression method on the sensor trajectory data and slicing on the MSA to achieve an improved anonymization model with a reduced information loss rate compared with earlier approaches.

The rest of the paper organizes as follows: Section II introduces the related works with their benefits, and Section III defines the basic definitions and notations incorporated. Section IV describes the procedure involved in Reciprocal Bucketization as the efficient anonymization model. The experimental results and comparative analysis with the existing approach are given in Section V. Finally, Section VI concludes the proposed approach.

## II. Literature Review

In this section presents the advantages of Smart Health Care, the recent research on the privacy preservation of trajectory data, and the multiple sensitive attributes addressing the various potential benefits and shortcomings.

Smart Health Care is a significant application of the IoT, where the various aspects of Health Care are implemented for ease of maintenance under the secured surveillance system. Vedaei *et al.* [16] presented COVID-SAFE, an automated health monitoring and surveillance system integrating IoT devices with Machine Learning algorithms. The primary aim is to improve the efficiency and accuracy of detecting and monitoring COVID-19 patients in urban areas. The depreciation of exposure to the coronavirus is more significant. Still, deployment and maintenance costs are high, and there needs to be a consideration for the security aspects of the data collected in the IoT environment.

IoT devices are highly vulnerable to attacks, and the medical data collected from those devices are highly significant and need highly secured privacy schemes and policies. Luo *et al.* [17] designed a Privacy Protector framework to defend against linkage attacks and secretly share data by adopting the Slepian-Wolf-coding-based secret sharing (SW-SSS). The distributed nature of the framework stores the patient data collected on the cloud server by assuring an efficient access control scheme. Privacy Protector ensures the security of the data collected, however securing the data in real-time is still challenging.

Komishani *et al.* [18] presents a Preserving personalized privacy in trajectory data publishing (PPTD) model for the trajectory data associated with the sensitive attribute of moving objects. The sensitive attribute generalization and trajectory data local suppression approach balance the data utility and privacy well. The linkages attacks such as identity linkage, attribute linkage, and similarity attacks are demonstrated on the anonymized data to the resistance of the data publishing. The PPTD has been implemented on the City80K and Metro100K datasets, and an extensive comparison is carried out with KCL. With less data loss and high privacy protection, the PPTD outstands in its performance aspects, but multiple sensitive attributes are yet to address with efficient data utility.

Addressing the various linkage attacks, such as attribute, record linkage, and similarity attacks on the trajectory data,

TABLE I. SMART HOSPITAL SENSOR TRAJECTORY DATASET WITH MULTIPLE SENSITIVE ATTRIBUTES ($P_{TB}$).

| PId | Trajectory | Multiple Sensitive Attributes | | | |
|---|---|---|---|---|---|
| | | Symptom | Diagnostic Method | Disease | Treatment |
| 1 | $x1 \rightarrow d2 \rightarrow z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | Abdominal pain | X-ray | Abdominal Cancer | Chemotherapy |
| 2 | $k6 \rightarrow n7 \rightarrow p8$ | Weight loss | Antibody Test | HIV | Medication |
| 3 | $z3 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | Eating disorders | Body mass index (BMI) | Obesity | Nutrition control |
| 1 | $x1 \rightarrow d2 \rightarrow n5 \rightarrow k6 \rightarrow p9$ | Fever | Molecular diagnostic methods | Cholera | Antibiotic |
| 4 | $x1 \rightarrow d2 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | Infection | ELISA Test | HIV | ART |
| 5 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow n7$ | Diarrhea | RT-PCR Tests | Dengue | Antibiotic |
| 6 | $x1 \rightarrow z3 \rightarrow n7 \rightarrow p8$ | Shortness of breath | FeNO test | Asthma | Medication |
| 7 | $x1 \rightarrow k6 \rightarrow n7 \rightarrow p8$ | Fever | RITD Tests | Influenza | Medicine |
| 8 | $n5 \rightarrow k6 \rightarrow p9$ | Weight loss | MRI Scan | Lung Cancer | Radiation Therapy |
| 9 | $d2 \rightarrow n5 \rightarrow n7 \rightarrow p9$ | Chest tightness | Methacholine challenge tests | Inflammation | Medication |
| 10 | $p4 \rightarrow n7 \rightarrow p8$ | Pain or discomfort | Biopsy Test | Skin cancer | Radiation Therapy |
| 11 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | Abdominal pain | Ultrasound | Dyspepsia | Antibiotic |

Yao *et al.* [19] have designed an anonymous technique called Enhanced l-diversity Data Privacy Preservation for publishing trajectory data (EDPP). EDPP defends against the background knowledge of the trajectory data to predict the sensitive attributes by identifying critical spatial-temporal sequences that cause privacy leakage. The method adopts perturbation and enhanced l - diversity with well-defined privacy constraints to ensure more excellent data utility. However, the approach must be extended for a greater trajectory length with indexing and multiple sensitive attributes.

Adding Noise to the trajectory data to ensure privacy through a vector-based grid environment is a new effort by Tojiboev *et al.* [20] Adding Noise before data publishing results in low complexity and greater privacy, but data handling and rebuilding the original data is challenging. Differential privacy is a modern, robust privacy preservation approach that implements a query mechanism to minimize privacy loss. Added Noise on the trajectory data and implementing differential privacy as the privacy protection model poses a challenge in data utility. They are considering the sensitive attribute label and adopting Generative Adversarial Network (GAN), an effective privacy preservation model implemented by Yao [21] The GAN ensures balanced data privacy and data utility. However, the computational speed and addressing the multi-source trajectory data and sensitive attributes are challenges for Differential Privacy.

Wen *et al.* [22] proposed a dynamic privacy level model by defining the relationship between privacy requirements and location features. The optimal differential privacy model is designed on the trajectory data by filtering the template trajectory with the semantic similarities on the trajectory as the constraint. The model publishes the data on randomization of the locations on the user trajectory data to balance privacy and data utility. However, adopting differential privacy on the multiple sensitive attributes and the dynamic trajectories is challenging.

Kanwal *et al.* [23] present (p, l)- Angelization for publishing 1:M, an individual with multiple records resisting a correlation attack. The Angelization method eliminates explicit Identifier by splitting the table into quasi-identifiers and multiple sensitive attributes. Quasi-attribute generalization and multiple sensitive attribute weight and dependency are computed for anonymization through (p, l)- Angelization. The algorithm performs well under static datasets without republications, but republication is an issue with data utility for the dynamic dataset.

The privacy preservation methods like generalization and bucketization pose a challenge to data utility. To overcome these issues slicing method is proposed by Li *et al.* [24] The slicing method can be applicable horizontally and vertically on the given data records; the membership disclosure attack is primarily defended. The significant advantage of slicing is to handle the high dimensional data with minimized data loss on the complete records. Overlapping slicing ensures high privacy by duplicating an attribute into more than one column. However, the utilization of the anonymized data still needs to be improved.

Heap Bucketization-anonymity (HBA) model is proposed by *et al.* [25], where the method develops an anonymization approach for quasi-identifiers and the sensitive attribute. HBA anatomizes the complete records to anonymize the sensitive attributes using slicing and Heap Bucketization of the quasi-identifier using k-anonymity and slicing. The KL - Divergence is used for validation in terms of utility and privacy by defending against background knowledge attacks, quasi-identifier attacks, membership attacks, and fingerprint correlation attacks. HBA results in less utility loss with greater privacy; HBA has yet to address the dynamic, unstructured data and semi-sensitive attributes.

Sensitive Label Privacy Preservation with Anatomization (SLPPA), a scheme for privacy preservation, is designed by Yao *et al.* [26] to address the various background knowledge attacks. The SLPPA adopts two phases in implementation, i.e., Table Division and Group Division. The entropy and mean-square contingency coefficient is computed for anonymizing by adding uncertainty during table division. The group division is performed by adopting the privacy constraints and ensuring no overlapping groups in the published data. SLPPA enhances data utility by defending against background knowledge attacks. However, dynamic data anonymization is yet challenging.

## III. PROBLEM DESCRIPTION AND BASIC NOTATIONS

### A. Problem Definition

The patient's dataset ($P_{TB}$), combines trajectory data ($P_{TB}^T$) and multiple sensitive attributes ($P_{TB}^{SA}$). Anonymize the dataset so that the Adversary with the prior knowledge fails to decode the individual identity through the trajectory data or the MSA. The anonymization approach must ensure the defense

mechanism against linkages attacks with optimal equilibrium between privacy and data utility.

### B. Notations

A patient's dataset consists of patient data records, where each record allows a unique patient identifier, sensory trajectory data, and a set of sensitive values. The patient sensor trajectory data with the MSA can be represented as:

$$P_{TB} = \{(PId_1, T_1, SA_1), \dots\dots(PId_i, T_i, SA_i)\} \quad (1)$$

Where, PId is unique identifier for a patient in $P_{TB}$. $T_i$ is an Sensor Trajectory data for a patient possessing the MSA as $SA_i$. $T_i$ is a sequence of data collected for the sensor's at a particular timestamp for a user i and it given as follows:

$$T_i = \{(sen_1, t_1)^i, (sen_2, t_2)^i, \dots(sen_n, t_n)^i\} \quad (2)$$

We compute $|T_i|$ as the number of sensors moving points for a patient 'i'. For an Example $|T_4| = 5$, i.e., five sensors are recording the data in sequence for patient with ID 4, concerning Table I.

Joinable trajectories are formed if there exists a sub-trajectory for the given trajectory. Let's consider $T_K = \{t_1^K, t_2^K, t_3^K, \dots, t_n^K\}$ as the sensor trajectory data and $T_S = \{t_1^s, t_2^s, t_3^s, \dots, t_m^s\}$ as the sub trajectory of user K then $|T_S|$ is the sub trajectory satisfying the condition n¡=m and $|T_S|$ must be subset of $|T_K|$. The trajectories can be merged using union operation.

### C. Adversary Model

The Adversary could breach the patient's privacy in two significant ways (1) the Adversary's Prior Knowledge of Trajectory Data and (2) the Adversary's goal to hit the victim by knowing a sensitive values from MSA. Let's consider that Avok is one of the patients in the smart hospital, and his details are recorded in Table I. Consider 'A' as the Adversary, which could be a data collector itself and aims to find the sensitive attributes of the targeted victim. With the prior knowledge about Avok, the adversary 'A' can perform the following attacks[27]:

Identity linkage attack: If adversary 'A' knows about Avok's sensor trajectory, i.e., sensor 'd' measures the oxygen level at timestamp 2 and sensor 'p' measure the blood pressure at timestamp 4, respectively, then A can claim that the record T1 belongs to Avok. The adversary can declare with 100% confidence that the record belongs to Avok and access the sensitive attributes because d2, p4 is the only sub-trajectory that belongs to the T1 record.

Attribute linkage attack: If adversary 'A' knows about Avok's partial sensor trajectory, i.e., sensor 'k' measures the heart rate at timestamp 6 and sensor 'n' measures the body temperature at timestamp 7, respectively. The table contains three records, T2, T4, and T5, with the sub-trajectory k6, n7 can be identified by A. Adversary 'A' can forecast that Avok has HIV disease with 67% confidence because out of three

records under multiple sensitive attributes, two records have the disease as HIV.

Correlated-records linkage attack: From the given Table, Avok has multiple records. If adversary 'A' knows the number of records Avok has in the dataset and also knowledge about Avok's sensor trajectory x1, k6. This makes the adversary predict with 100% confidence with his prior correlated knowledge that Avok has Asthma plus dengue could be specifically on the Sensor trajectory data. Similarly, if the Adversary 'A' knows the two sensitive attributes out of four, like the Symptom: Abdominal pain and Treatment: Chemotherapy, the adversary could correlate them to say confidently that the patient has Abdominal Cancer.

### D. Privacy Requirement

The goal of anonymizing the patient dataset from adversaries requires adopting the following privacy requirements in our proposed approach [28].

*1) Bucketization:* The patient's dataset $(P_{TB})$, which is a combination of sensor trajectory data $(P_{TB}^T)$ and multiple sensitive attributes $(P_{TB}^{SA})$, is partitioned into 'n' buckets with a constraint of a maximum bucket size of three patient records. On partitioned, each bucket is named with Bucket ID. Further bucketization helps in carrying the suppression and slicing parallelly.

*2) Suppression:* Suppression is the method of eliminating the critical sensor trajectory point from the patient trajectory data. A sensor trajectory point is critical if and only if the trajectory point fails to satisfy the K-1 threshold defined and the suppression metric. The critical point is eliminated only with the corresponding trajectory to enhance the data utility by minimizing the data loss, which could express as Local Suppression [29].

*3) Slicing:* It's a method to anonymize the data using a partition. Partition could be carried vertically or horizontally to get the randomly permutated sliced table. The attribute partition is carried out by slicing Table III vertically into two slices, say (Disease, Symptom) and (Diagnostic Method, Treatment) [30].

### E. Utility Metrics

To maintain the trade-off between the patient's privacy and data utility in the published anonymized dataset. It's required to define the Reciprocal Bucketization to ensure less data loss through cumulated sensitive attribute representation under the same bucket. The suppression metric measures the suppression score of every sensor trajectory point and helps find the critical trajectory point to remove. The computation of the suppression score of one sensor point from the critical trajectory is as follows.

$$SupressionScore = \frac{NT\_B_i\_CP}{NT\_B_i} \quad (3)$$

Where, $NT\_B_i\_CP$ : Number of Trajectories in Bucket with the Critical Point and $NT\_B_i$ : Number of Trajectories in the bucket

On the computation of the suppression score for both the sensor points from the trajectory, one point is selected as critical by finding the maximum among them. If both the points have the same suppression score, the leftmost sensor trajectory point is declared the critical point to remove [31].

## IV. PROPOSED SYSTEM

The reciprocal Bucketization anonymity model proposes the design workflow as shown in the Fig. 2 and algorithm to prevent the privacy leakage of patients' data from the three linkage attacks. The proposed model consists of the following four phases,
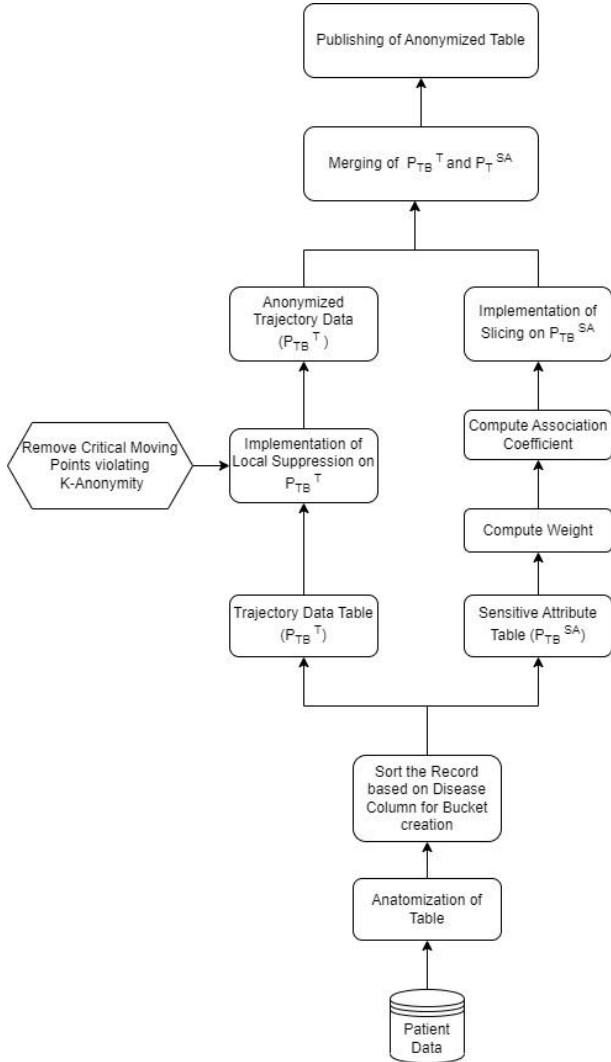


Fig. 2. Workflow of reciprocal bucketization model.

1) Bucketization and Partition.
2) Local Suppression on Sensor Trajectory Data ($P_{TB}^T$).
3) Slicing of MSA ($P_{TB}^{SA}$).
4) Publishing Anonymized Dataset on Merging $P_{TB}^T$ and $P_{TB}^{SA}$

The primary goal of the proposed model is to ensure less data loss with high privacy preservation. The detailed procedure for the above four phases is in the following sub-sections.

---

**Algorithm 1** Reciprocal Bucketization

**INPUT:** Patient Data $P_{TB}$, Bucket_Size
**OUTPUT:** Anonymised Patient Data $P_{TB}^A$.

1: anatomize (Patient Data $P_{TB}$)
2: Sort_Disease_Value(Patient Data $P_{TB}$)
3: $Bucket\_ID = 0, count = 0, i = 1$
4: **for** $k \leftarrow 1$ to $len(P_{TB})$ **do**
5:     **if** $count \leq Bucket\_Size$ **then**
6:         $Bucket\_ID = i$
7:         $count = count + 1$
8:     **else**
9:         $count = 0$
10:         $i = i + 1$
11:     **end if**
12: **end for**
13: $SplitP_{TB} = P_{TB}^T withBucket\_ID, P_{TB}^{SA} withBucket\_ID$
14: $P_{TB}^{TA} = $ Anonymise_Trajectory($P_{TB}^T$)
15: $P_{TB}^{SAA} = $ Anonymise_Sensitive_Attribute($P_{TB}^{SA}$)
16: $P_{TB}^A = $ merge($P_{TB}^{TA}, P_{TB}^{SAA}$)
17: return $P_{TB}^A$

---

### A. Bucketization and Partition

The primary goal of this step is to create the Bucket based on the user input bucket size, then partition Table II into two table's Table III and Table V, having sensor trajectory and MSA, respectively. Before the bucket creation, the entire Table I is sorted according to the Disease column from the MSA to get Table II and a new column, Bucket ID, where each patient's record is added with the bucket number sequentially corresponding to the bucket size. I.e., if the user input bucket size is 3, then the first three records from Table II are allocated with bucket number 1, and so on for the rest of the buckets. Further, Table II is divided into two Tables to carry Suppression on the Trajectory data and Slicing on MSA to achieve the anonymization effectively.

### B. Local Suppression on Sensor Trajectory Data ($P_{TB}^T$)

On the sensory trajectory data $P_{TB}^T$, the Local Suppression method is implemented to remove the critical trajectory points from the patient's PTB dataset and generate the anonymized trajectory dataset. The algorithms depict the steps involved in the trajectory suppression, and the procedure is as follows:

TABLE III. SENSOR TRAJECTORY DATA ($P_{TB}^T$)

| PId | Trajectory | Bucket ID |
|---|---|---|
| 1 | $x1 \rightarrow d2 \rightarrow z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | 1 |
| 6 | $x1 \rightarrow z3 \rightarrow n7 \rightarrow p8$ | 1 |
| 1 | $x1 \rightarrow d2 \rightarrow n5 \rightarrow k6 \rightarrow p9$ | 1 |
| 5 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow n7$ | 2 |
| 11 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | 2 |
| 2 | $k6 \rightarrow n7 \rightarrow p8$ | 2 |
| 4 | $x1 \rightarrow d2 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | 3 |
| 9 | $d2 \rightarrow n5 \rightarrow n7 \rightarrow p9$ | 3 |
| 7 | $x1 \rightarrow k6 \rightarrow n7 \rightarrow p8$ | 3 |
| 8 | $n5 \rightarrow k6 \rightarrow p9$ | 4 |
| 3 | $z3 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | 4 |
| 10 | $p4 \rightarrow n7 \rightarrow p8$ | 4 |

TABLE II. SORTING DISEASE COLUMN AND BUCKET ID ASSIGNMENT

| PId | Trajectory | Multiple Sensitive Attributes | | | | |
|---|---|---|---|---|---|---|
| | | Symptom | Diagnostic Method | Disease | Treatment | Bucket ID |
| 1 | $x1 \rightarrow d2 \rightarrow z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | Abdominal pain | X-ray | Abdominal Cancer | Chemotherapy | 1 |
| 6 | $x1 \rightarrow z3 \rightarrow n7 \rightarrow p8$ | Shortness of breath | FeNO test | Asthma | Medication | 1 |
| 1 | $x1 \rightarrow d2 \rightarrow n5 \rightarrow k6 \rightarrow p9$ | Fever | Molecular diagnostic methods | Cholera | Antibiotic | 1 |
| 5 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow n7$ | Diarrhea | RT-PCR Tests | Dengue | Antibiotic | 2 |
| 11 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | Abdominal pain | Ultrasound | Dyspepsia | Antibiotic | 2 |
| 2 | $k6 \rightarrow n7 \rightarrow p8$ | Weight loss | Antibody Test | HIV | Medication | 2 |
| 4 | $x1 \rightarrow d2 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | Infection | ELISA Test | HIV | ART | 3 |
| 9 | $d2 \rightarrow n5 \rightarrow n7 \rightarrow p9$ | Chest tightness | Methacholine challenge tests | Inflammation | Medication | 3 |
| 7 | $x1 \rightarrow k6 \rightarrow n7 \rightarrow p8$ | Fever | RITD Tests | Influenza | Medicine | 3 |
| 8 | $n5 \rightarrow k6 \rightarrow p9$ | Weight loss | MRI Scan | Lung Cancer | Radiation Therapy | 4 |
| 3 | $z3 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | Eating disorders | Body mass index (BMI) | Obesity | Nutrition control | 4 |
| 10 | $p4 \rightarrow n7 \rightarrow p8$ | Pain or discomfort | Biopsy Test | Skin cancer | Radiation Therapy | 4 |

The algorithm accepts $P_{TB}^{T}$ sensor trajectory data with Bucket IDs, adversary prior knowledge delta, and threshold K to produce anonymized sensor trajectory data as the output. In the suppression procedure, the primary step is to group all the trajectory data into two groups: one with a similar Bucket ID for which the anonymization is to be carried and the other group with all remaining Bucket IDs. Split the trajectory data into two sets, say set X consists of trajectory data with '1' as the Bucket ID, and set Y consists of the records with remaining Bucket ids. Now from set X, find all the sub-trajectories of length 1 and verify that those trajectories appeared with a minimum of K-1 times in set Y. If any of the sub-trajectories fails to leave traces of K-1 times, those trajectories are eliminated; else, trajectories are preserved as same.

TABLE IV. ANONYMIZED TRAJECTORY DATA $(P_{TB}^{TA})$

| PId | Trajectory | Bucket ID |
|---|---|---|
| 1 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | 1 |
| 6 | $z3 \rightarrow n7 \rightarrow p8$ | 1 |
| 1 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow p9$ | 1 |
| 5 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow n7$ | 2 |
| 11 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | 2 |
| 2 | $k6 \rightarrow n7 \rightarrow p8$ | 2 |
| 4 | $d2 \rightarrow k6 \rightarrow p9$ | 3 |
| 9 | $d2 \rightarrow n5 \rightarrow p9$ | 3 |
| 7 | $k6 \rightarrow n7 \rightarrow p8$ | 3 |
| 8 | $n5 \rightarrow k6 \rightarrow p9$ | 4 |
| 3 | $z3 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | 4 |
| 10 | $p4 \rightarrow p8$ | 4 |

Similarly, find all the sub-trajectories of length 2 from set X and validate with set Y for the minimum traces to K-1 times, failing to eliminate a trajectory point that does not satisfy the condition by computing the suppression score. Else keep the trajectory the same in the record. Repeat the procedure to calculate the critical trajectory point till the length of the sub-trajectory equals the adversary prior knowledge length delta. The complete process has to iterate for all the Bucket ID's for the PTBT to generate an anonymized sensor trajectory record PTBTA, as shown in Table IV.

Consider an example of the sensor trajectory $x1 \rightarrow d2 \rightarrow z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ , which belongs to Bucket ID 1. The sensor trajectory is validated by computing the suppression score for critical trajectory points of the length 1 and 2, respectively, and found that the trajectory points x1 and d2 are critical to eliminating. On the complete trajectory iteration, the anonymized trajectory is $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$.

---

**Algorithm 2** Anonymise_Trajectory($P_{TB}^{T}$)

---

**INPUT:** Sensor-Trajectory data $P_{TB}^{T}$ with Bucket_ID, $A's$ prior knowledge $\partial$ with maximum length $\rho$, K threshold
**OUTPUT:** Anonymized Sensor Trajectory data $P_{TB}^{TA}$.

1: **Scan** Sensor Trajectory Table $P_{TB}^{T}$
2: let $S_A$ = {set of all distinct sensitive values under the same Bucket_ID}
3: **for** each $sa \in S_A$ **do**
4:     $i = 1$, $\text{Cr}_p = \emptyset$, $\text{Dr}_i = \emptyset$
5:     $\text{P} = \{P_{TBr}^{T} \mid P_{TBr}^{T} \in P_{TB}^{T} \wedge P_{TBr}^{T}(sa) = sa\}$
6:     $\text{Q} = P_{TBr}^{T} - \{\text{P}\}$
7:     **for** each $P_{TBr}^{T} \in \text{P}$ **do**
8:         $\text{Cr}_p = \{\tau_r \mid \tau_r \subseteq P_{TBr}^{T} \wedge \mid \tau_r \mid = 1\}$
9:         **for** each $\tau_r \in \text{Cr}_p$ **do**
10:             **if** $(\mid \tau_r \in P_{TBr}^{T} \mid_{\forall P_{TBr}^{T} \in \text{Q}} \geq \text{K})$ **then**
11:                 $\text{Dr}_i = \text{Dr}_i \cup \tau_r$
12:             **else**
13:                 remove $\tau_r$ from $P_{TBr}^{T} \in \text{P}$
14:             **end if**
15:         **end for**
16:     **end for**
17:     **while** $(i + 1 \leq \rho)$ **do**
18:         **for** each $\tau_r \in \text{Dr}_i$ join with successive $\tau_{r+i}$ in $\text{Dr}_i$ **do**
19:             **if** $(\mid \tau_r \cup \tau_{r+i} \in P_{TBr}^{T} \mid_{\forall P_{TBr}^{T} \in \text{Q}} \geq \text{k})$ **then**
20:                 $\text{Dr}_{i+1} = \text{Dr}_{i+1} \cup \{\tau_r \cup \tau_{r+i}\}$
21:             **else**
22:                 $\text{tr} = \tau_r \cup \tau_{r+i}$
23:                 remove $\chi(t)$ from $P_{TBr}^{T} \in \text{P}$
24:             **end if**
25:         **end for**
26:         $i = i + 1$
27:     **end while**
28: **end for**

---

### C. Slicing of Multiple Sensitive Attributes $(P_{TB}^{SA})$

Anonymizing the MSA $P_{TB}^{SA}$ obtained after Bucketization and Partition implements the following slicing steps. In the slicing procedure, we first measure the weight of each sensitive attribute through the concept of entropy. Entropy refers to the average value of the information in each message gained. The sensitive attribute with higher entropy measures results as the qualitative information container. The entropy is measured with the following formula:

$$W_{SA} = - \sum_{j=1}^{D_{SA}} p(S_{vi}) \log(p(S_{vi})) \quad (4)$$

TABLE V. MULTIPLE SENSITIVE ATTRIBUTE ($P_{TB}^{SA}$)

| PId | Symptom | Diagnostic Method | Disease | Treatment | Bucket ID |
|---|---|---|---|---|---|
| 1 | Abdominal pain | X-ray | Abdominal Cancer | Chemotherapy | 1 |
| 6 | Shortness of breath | FeNO test | Asthma | Medication | 1 |
| 1 | Fever | Molecular diagnostic methods | Cholera | Antibiotic | 1 |
| 5 | Diarrhea | RT-PCR Tests | Dengue | Antibiotic | 2 |
| 11 | Abdominal pain | Ultrasound | Dyspepsia | Antibiotic | 2 |
| 2 | Weight loss | Antibody Test | HIV | Medication | 2 |
| 4 | Infection | ELISA Test | HIV | ART | 3 |
| 9 | Chest tightness | Methacholine challenge tests | Inflammation | Medication | 3 |
| 7 | Fever | RITD Tests | Influenza | Medicine | 3 |
| 8 | Weight loss | MRI Scan | Lung Cancer | Radiation Therapy | 4 |
| 3 | Eating disorders | Body mass index (BMI) | Obesity | Nutrition control | 4 |
| 10 | Pain or discomfort | Biopsy Test | Skin cancer | Radiation Therapy | 4 |

Where SA is Sensitive Attribute, $\{S_{v1}, S_{v2}, S_{v3}, ...., S_{vi}\}$ set of possible values under SA. $p(S_{vi})$ possibility that $S_{vi}$ is considered and $D_{SA}$ number of distinct sensitve attributes.

The slicing procedure continues to find the association after computing the weight of each sensitive attribute using entropy. We adopt Pearson's Contingency Coefficient in the association computation to measure the association between two sensitive attributes. The analysis obeys the following formula:

$$\phi^2(SA_1, SA_2) = \frac{\sum_{i=1}^{n1} \sum_{j=1}^{n2} \frac{(p(SA_{ij}) - p(SA_i)p(SA_j))^2}{p(SA_i)p(SA_j)}}{min\{n1, n2\} - 1}) \quad (5)$$

where, n1 and n2 are the total number of distinct values of $SA_1$ and $SA_2$ respectively. $p(SA_{ij})$ represents the chance from $SA_{ij}$. $p(SA_i)$ and $p(SA_j)$ are the boundary totals, where $p(SA_i) = \sum_{j=1}^{n2}$ and $p(SA_j) = \sum_{i=1}^{n1}$

On the computation of the weights of each sensitive attribute and the association among them, we combine the four columns of the MSA to form two columns for generating the sliced table. The sensitive attribute with the higher weights is selected as the slice's first column to avoid the adversaries' ease of data access. Then we choose the second column of the slice, with the maximum average association coefficient, with the other column. Doing so maximizes the association between the sensitive attributes in the same silce, and the coefficient for attributes in the different slices is minimized. After completing the above procedure, Table V with sensitive attributes like Symptoms, Diagnostic Method, Disease, and Treatment generates Table VI. Table VI has two slices on the attributes (Disease, Symptom) and (Diagnostic Method, and Treatment) to efficiently anonymize the data.

### D. Publishing Anonymized Dataset on Merging $P_{TB}^{TA}$ and $P_{TB}^{SAA}$

The anonymized sensor trajectory data and MSA data are merged to get the Table to give the input to the Reciprocal Bucketization. In this procedure, the sensor trajectory data remains untouched. In return to the Bucket ID, the MSA is combined to form three tuples having the three patients' data records under the same Bucket ID, which can reduce the adversary's confidence in identifying the individual through sensor trajectory data or the multiple sensitive attributes. In merging the corresponding records, the Patient ID is the

---

**Algorithm 3** Anonymise_Sensitive_Attribute($P_{TB}^{SA}$)

**INPUT:** Multiple Sensitive Attributes $P_{TB}^{SA}$ with Bucket_ID
**OUTPUT:** Anonymized Multiple Sensitive Attributes $P_{TB}^{SAA}$.

1: **Scan** Multiple Sensitive Attributes Table $P_{TB}^{SA}$
2: **for** each attribute in $P_{TB}^{SA}$ **do**
3:     compute $W_{SA}$ store in $W_{SA}Array$
4: **end for**
5: $SA\_Max_1 = First\_Max(W_{SA}Array)$
6: $SA\_Max_2 = Second\_Max(W_{SA}Array)$
7: $Slice\_01 = (SA\_Max_1)$
8: $Slice\_02 = (SA\_Max_2)$
9: Compute Pearson Contingency Coefficient of MSA excluding $(SA\_Max_1) and (SA\_Max_2)$
10: compute $\phi^2(SA_1, SA_2)$ store in $PCC_{SA}Array$
11: $PCC\_Max_1 = First_M ax(PCC_{SA}Array)$
12: $PCC\_Max_2 = Second_M ax(PCC_{SA}Array)$
13: $Slice\_01 = (SA\_Max_1, PCC\_Max_1)$
14: $Slice\_02 = (SA\_Max_2, PCC\_Max_2)$
15: return($Slice\_01, Slice\_02$)

---

referral point. Finally, the anonymized data is published by sorting the records according to the patient ID as represented in the Table VII.

The RB – Anonymization model generates the anonymized dataset to be published, as represented in Table VII, and it is resistant to Identity, Attribute, and Correlated linkage attacks. To our knowledge, RB -Anonymization model is the first approach to anonymize the trajectory data with the MSA. If adversary prior expertise with the sensor trajectory length delta = 2 as d2,p4, then Adversary 'A' can perform all three linkage attacks on Table I as discussed in Section III. However, the Adversary fails to identify Avok's record from Table VII because the not even one record with d2,p4. Only by finding the sensory trajectory point d2, the Adversary could declare the identification with less than 37% confidence. The MSA couldn't be reached with the same Adversary's prior knowledge.

Similarly, knowing the partial sensor trajectory k6, n7 with the significant MSA as the disease HIV could be predicted in Table I. But from Table VII even though the k6, n7 trajectory is repeated more than one time, the adversary failed to identify its

TABLE VI. SLICED MULTIPLE SENSITIVE ATTRIBUTE ($P_{TB}^{SAA}$)

| PId | (Disease, Symptom) | (Diagnostic Method , Treatment) | Bucket ID |
|---|---|---|---|
| 1 | (Abdominal Cancer , Abdominal pain) | (X-ray , Chemotherapy) | 1 |
| 6 | (Cholera , Fever) | (Molecular diagnostic methods , Antibiotic) | 1 |
| 1 | (HIV , Weight loss) | (Antibody Test , Medication) | 1 |
| 5 | (Obesity , Eating disorders) | (Body mass index (BMI) , Nutrition control) | 2 |
| 11 | (HIV , Infection) | (ELISA Test , ART) | 2 |
| 2 | (Dengue , Diarrhea) | (RT-PCR Tests , Antibiotic) | 2 |
| 4 | (Asthma , Shortness of breath) | (FeNO test , Medication) | 3 |
| 9 | (Influenza , Fever) | (RITD Tests , Medicine) | 3 |
| 7 | (Lung Cancer , Weight loss) | (MRI Scan , Radiation Therapy) | 3 |
| 8 | (Inflammation , Chest Tightness) | (Methacholine challenge tests , Medication) | 4 |
| 3 | (Skin cancer , Pain or discomfort) | (Biopsy Test , Radiation Therapy) | 4 |
| 10 | (Dyspepsia , Abdominal pain) | (Ultrasound , Antibiotic) | 4 |

TABLE VII. MERGING AND RECIPROCAL BUCKETIZATION $P_{TB}^{A}$

| PId | Trajectory | (Disease, Symptom) | (Diagnostic Method , Treatment) |
|---|---|---|---|
| 1 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | (Abdominal Cancer , Abdominal pain)<br>(Cholera , Fever)<br>(HIV , Weight loss) | (X-ray , Chemotherapy)<br>(Molecular diagnostic methods , Antibiotic)<br>(Antibody Test , Medication) |
| 2 | $k6 \rightarrow n7 \rightarrow p8$ | (Obesity , Eating disorders)<br>(HIV , Infection)<br>(Dengue , Diarrhea) | (Body mass index (BMI) , Nutrition control)<br>(ELISA Test , ART)<br>(RT-PCR Tests , Antibiotic) |
| 3 | $z3 \rightarrow k6 \rightarrow n7 \rightarrow p9$ | (Inflammation , Chest Tightness)<br>(Skin cancer , Pain or discomfort)<br>(Dyspepsia , Abdominal pain) | (Methacholine challenge tests , Medication)<br>(Biopsy Test , Radiation Therapy)<br>(Ultrasound , Antibiotic) |
| 1 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow p9$ | (Abdominal Cancer , Abdominal pain)<br>(Cholera , Fever)<br>(HIV , Weight loss) | (X-ray , Chemotherapy)<br>(Molecular diagnostic methods , Antibiotic)<br>(Antibody Test , Medication) |
| 4 | $d2 \rightarrow k6 \rightarrow p9$ | (Asthma , Shortness of breath)<br>(Influenza , Fever)<br>(Lung Cancer , Weight loss) | (FeNO test , Medication)<br>(RITD Tests , Medicine)<br>(MRI Scan , Radiation Therapy) |
| 5 | $d2 \rightarrow n5 \rightarrow k6 \rightarrow n7$ | (Obesity , Eating disorders)<br>(HIV , Infection)<br>(Dengue , Diarrhea) | (Body mass index (BMI) , Nutrition control)<br>(ELISA Test , ART)<br>(RT-PCR Tests , Antibiotic) |
| 6 | $z3 \rightarrow n7 \rightarrow p8$ | (Abdominal Cancer , Abdominal pain)<br>(Cholera , Fever)<br>(HIV , Weight loss) | (X-ray , Chemotherapy)<br>(Molecular diagnostic methods , Antibiotic)<br>(Antibody Test , Medication) |
| 7 | $k6 \rightarrow n7 \rightarrow p8$ | (Asthma , Shortness of breath)<br>(Influenza , Fever)<br>(Lung Cancer , Weight loss) | (FeNO test , Medication)<br>(RITD Tests , Medicine)<br>(MRI Scan , Radiation Therapy) |
| 8 | $n5 \rightarrow k6 \rightarrow p9$ | (Inflammation , Chest Tightness)<br>(Skin cancer , Pain or discomfort)<br>(Dyspepsia , Abdominal pain) | (Methacholine challenge tests , Medication)<br>(Biopsy Test , Radiation Therapy)<br>(Ultrasound , Antibiotic) |
| 9 | $d2 \rightarrow n5 \rightarrow p9$ | (Asthma , Shortness of breath)<br>(Influenza , Fever)<br>(Lung Cancer , Weight loss) | (FeNO test , Medication)<br>(RITD Tests , Medicine)<br>(MRI Scan , Radiation Therapy) |
| 10 | $p4 \rightarrow p8$ | (Inflammation , Chest Tightness)<br>(Skin cancer , Pain or discomfort)<br>(Dyspepsia , Abdominal pain) | (Methacholine challenge tests , Medication)<br>(Biopsy Test , Radiation Therapy)<br>(Ultrasound , Antibiotic) |
| 11 | $z3 \rightarrow p4 \rightarrow k6 \rightarrow p8$ | (Obesity , Eating disorders)<br>(HIV , Infection)<br>(Dengue , Diarrhea) | (Body mass index (BMI) , Nutrition control)<br>(ELISA Test , ART)<br>(RT-PCR Tests , Antibiotic) |

MSA due to the bucketization where three tuples are confusing to decide, and the probability of finding is less than 0.2, which is slightly negligible.

Correlated linkage attacks concerning multiple sensory trajectories with the same patient ID and establishing the relations among the sensitive attributes are well addressed by RB – Anonymization. The significant bucketization process brings down the confidence level of the adversary in the correlated linkage attack to less than 30% as the outcome of Bucketization and combining the records under the same bucket ID.

## V. RESULTS AND DISCUSSION

The proposed approach uses the Windows 10 operating system, i5 processor with a minimum of 4GB RAM and 256 GB SSD. The algorithm implementation uses the Python 3 version of the Synthetic dataset. The number of instances considered is more excellent than 10,000. The dataset's attributes are PId, Sensor Trajectory, Disease, Symptom, Diagnostic Method, and Treatment. Disease, Symptom, Diagnostic Method, and Treatment categorizes as MSA. The proposed approach's performance evaluates in terms of information loss in anonymized sensor trajectory data and the utility loss in MSA ($P_{TB}^{A}$).

### A. Sensor Trajectory Information Loss

Information loss occurs during the anonymization process because of the methods used, like generalization or suppression. Analysing the number of sensor trajectories information loss in the resultant anonymized dataset is quite significant. Eliminating the critical trajectory moving point results in
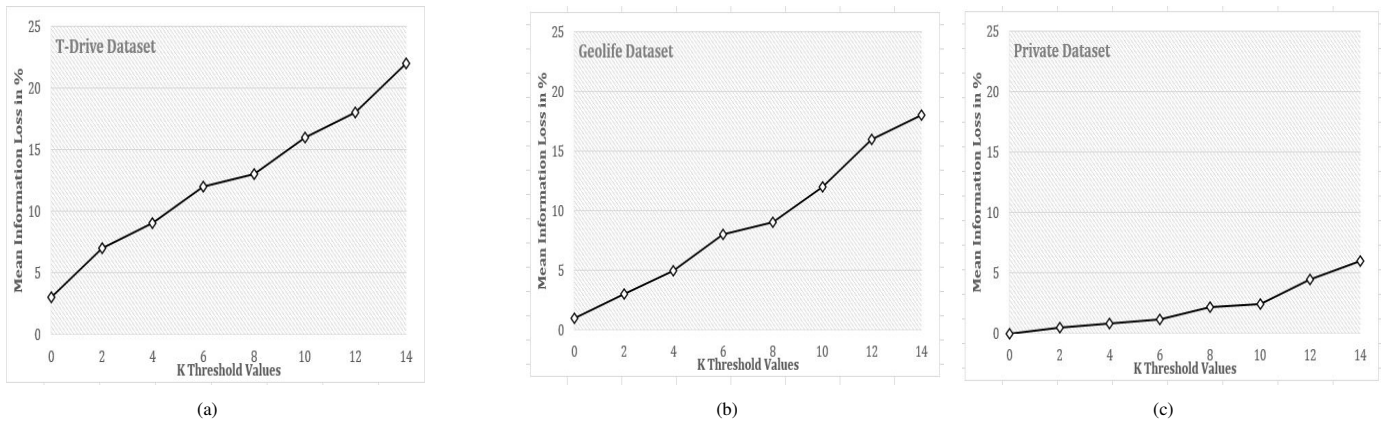
Fig. 3. Mean sensor trajectory information loss in $P_{TB}^{TA}$ with K threshold values. (a) T- Drive dataset. (b) Geolife dataset. (c) Private dataset.
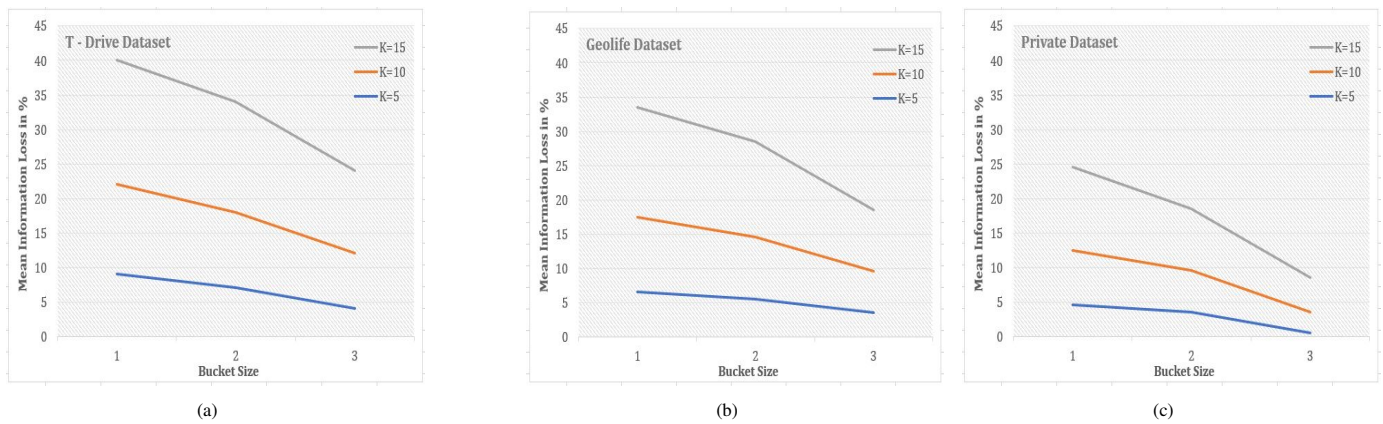


Fig. 4. Mean sensor trajectory information loss in $P_{TB}^{TA}$ with bucket size. (a) T- Drive dataset. (b) Geolife dataset. (c) Private dataset.

information loss during the procedure to satisfy the privacy requirements.

The sensor trajectory information loss is computed as follows between the $P_{TB}^{T}$ (Original Sensor Trajectory) and $P_{TB}^{TA}$ (Anonymized Sensor Trajectory).

Information Loss($I_L$) of each record is computed as:

$$I_L(P_{TBi}^{TA}) = \frac{|P_{TBi}^{T}| - |P_{TBi}^{TA}|}{|P_{TBi}^{T}|} \qquad (6)$$

Total Trajectory Information Loss is computed as:

$$I_L(P_{TB}^{TA}) = \sum_{i=1}^{|P_{TB}^{TA}|} I_L(P_{TBi}^{TA}) \qquad (7)$$

where, $P_{TBi}^{T}$ represents total number of trajectory points in $i^{th}$ record of $P_{TB}^{T}$ and $P_{TBi}^{TA}$ represents total number of trajectory points in $i^{th}$ of $P_{TB}^{TA}$.

Fig. 3 shows the mean sensor trajectory information loss in the seven sensor trajectories corresponding to the various K threshold values. The graph shows that the increase in K

values is directly proportional to the information loss on the anonymized data due to the random increase in the critical trajectory points, which failed to reach the privacy requirements. Hence the data publisher has to adopt the K value carefully to hold the moderate data loss and to preserve privacy.

Effect of Bucket Size: The RB anonymity model splits the given records in the user input bucket size, where we have taken as 3. The number of sensor trajectories that fall under the bucket is directly proportional to the bucket size. Anonymizing the trajectories while keeping the constraints of K threshold values and adversary prior knowledge varies on the bucket size. Fig. 4 represents the mean Information loss of the patient's sensor trajectory by keeping the adversary prior knowledge sigma to 2 and changing the K threshold values. The significant observation from the graph is that as the bucket's size increases, the information loss reduces due to the more trajectories falling into the bucket and the less elimination of trajectory critical point. It leads to quick access to the sensitive attributes hence the bucket size has to be chosen appropriately by providing equal significance to the sensor trajectory and the multiple sensitive attributes.

## B. Utility Loss in Multiple Sensitive Attributes

The utility loss in MSA measures probability distribution across the actual values and the RB – Anonymize data. We adopt the Kullback – Leibler divergence metric to estimate the probability distribution differences. The MSA $P_{TB}^{SA}$ table is considered as actual values for implementing the KL divergence as x1. x1(r) are the elements of the records for R, i.e., (r belongs to R). x2 is an estimated probability distribution considered after anonymization ($P_{TB}^{A}$). The KL divergence is given as follows:

$$KL_{div}(x1, x2) = \sum_{r \in R} x1(r) \log(\frac{x1(r)}{x2(r)}) \qquad (8)$$
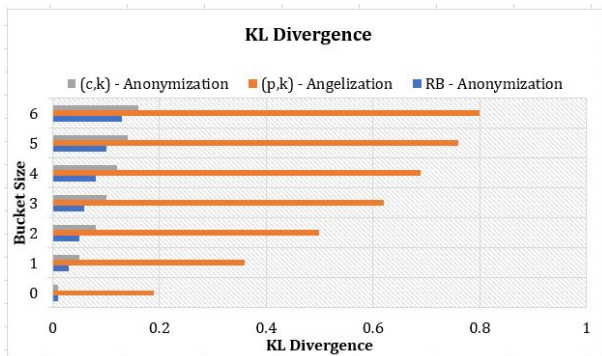


Fig. 5. KL Divergence of MSA.

Fig. 5 represents the utility loss comparison between the existing methods discussed in Section II with the proposed RB – Anonymization model. We can notify the significant differences, such as the (p,k) – angelization may assure high privacy. Still, the estimated probability distribution increases, resulting in a moderate utility loss. (c,k) – Anonymization and RB – Anonymization approach has a slight difference in the probability distribution with negligible utility loss; RB – Anonymization model outperforms consistently even if bucket size increases.

The analysis of the RB – Anonymization is validated by measuring the execution time and the privacy loss. Since no exact dataset has the sensor trajectory data and the multiple sensitive attributes, we have built our own synthetic dataset to evaluate the approach. The proposed model ensures its high adaptability towards real-time applications through its little execution time for anonymization with greater privacy. A moderate number of records and sensitive attributes are preferable since the execution time is proportional to both attributes. As the RB – Anonymization approach outperforms in terms of privacy, it also requires assuring the privacy loss on publishing the anonymized data. The privacy loss is negligible if there is less patient data exposure in terms of sensor trajectory and the MSA. The proposed approach brings new challenges to the adversaries through its stringent privacy preservation approach.

## VI. Conclusions

In this paper, we present a novel privacy preservation method considering the Smart Hospital data, consisting of sensor trajectory and the MSA. The proposed method outperforms in defending Identity, Attribute, and Correlated linkage attacks on data publishing. Our approach adopts a local suppression to anonymize the sensor trajectory and the slicing for MSA with a constraint on the bucket size. The proposed method outperforms on comparing the information loss of the sensor trajectory and MSA with (p,k) – angelization and (c,k) – anonymization approaches implemented on the real-time and synthetic dataset. As a future study, we are interested in addressing all the various linkage attacks and, as the primary concentration on the MSA, enabling our proposed method as the best practice.

## References

[1] G. Xu, "IoT-Assisted ECG Monitoring Framework with Secure Data Transmission for Health Care Applications," *IEEE Access*, vol. 8, pp. 74 586–74 594, 2020.

[2] Y. Shi, Z. Zhang, H. C. Chao, and B. Shen, "Data Privacy Protection Based on Micro Aggregation with Dynamic Sensitive Attribute Updating," *Sensors (Switzerland)*, vol. 18, no. 7, pp. 1–16, 2018.

[3] A. Majeed and S. Lee, "Anonymization Techniques for Privacy Preserving Data Publishing: A Comprehensive Survey," *IEEE Access*, vol. 9, pp. 8512–8545, 2021.

[4] N. Y. Philip, M. Razaak, J. Chang, M. S. Suchetha, M. Okane, and B. K. Pierscionek, "A Data Analytics Suite for Exploratory Predictive, and Visual Analysis of Type 2 Diabetes," *IEEE Access*, vol. 10, pp. 13 460–13 471, 2022.

[5] R. Khan, X. Tao, A. Anjum, H. Sajjad, S. U. R. Malik, A. Khan, and F. Amiri, "Privacy Preserving for Multiple Sensitive Attributes against Fingerprint Correlation Attack Satisfying c -Diversity," *Wireless Communications and Mobile Computing*, vol. 2020, pp. 1–18, 2020.

[6] N. L. and Raju, M. Seetaramanath, and P. S. Rao, "A Novel Dynamic KCi - Slice Publishing Prototype for Retaining Privacy and Utility of Multiple Sensitive Attributes," *International Journal of Information Technology and Computer Science*, vol. 11, no. 4, pp. 18–32, 2019.

[7] E. G. Komishani and M. Abadi, "A Generalization-Based Approach for Personalized Privacy Preservation in Trajectory Data Publishing," *In the Proceedings of 6th International Symposium on Telecommunications, IST 2012*, pp. 1129–1135, 2012.

[8] Y. Alotaibi, "A New Secured E-Government Efficiency Model for Sustainable Services Provision," *Journal of Information Security and Cybercrimes Research*, vol. 3, no. 1, pp. 75–96, 2020.

[9] A. Ye, Q. Zhang, Y. DIao, J. Zhang, H. Deng, and B. Cheng, "A Semantic-Based Approach for Privacy- Preserving in Trajectory Publishing," *IEEE Access*, vol. 8, pp. 184 965–184 975, 2020.

[10] F. Jin, W. Hua, M. Francia, P. Chao, M. Orowska, and X. Zhou, "A Survey and Experimental Study on Privacy-Preserving Trajectory Data Publishing," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 6, pp. 5577–5596, 2022.

[11] D. Hemkumar, S. Ravichandra, and D. V. Somayajulu, "Impact of Prior Knowledge on Privacy Leakage in Trajectory Data Publishing," *Engineering Science and Technology, an International Journal*, vol. 23, no. 6, pp. 1291–1300, 2020. [Online]. Available: https://doi.org/10.1016/j.jestch.2020.06.002

[12] J. N. Vanasiwala and N. R. Nanavati, "Multiple sensitive attributes based privacy preserving data publishing," *In the Proceedings of the 2nd International Conferenceon Computing Methodologies and Communication, ICCMC*, pp. 394–400, 2018.

[13] Y. Xiao and H. Li, "Privacy Preserving Data Publishing for Multiple Sensitive Attributes Based on Security Level," *Information (Switzerland)*, vol. 11, no. 3, p. 166, 2020.

[14] X. Liu and Y. Zhu, "Privacy and Utility Preserving Trajectory Data Publishing for Intelligent Transportation Systems," *IEEE Access*, vol. 8, pp. 176 454–176 466, 2020.

[15] J. Zhao, J. Mei, S. Matwin, Y. Su, and Y. Yang, "Risk-Aware Individual Trajectory Data Publishing with Differential Privacy," *IEEE Access*, vol. 9, pp. 7421–7438, 2021.

[16] S. S. Vedaei, A. Fotovvat, M. R. Mohebbian, G. M. Rahman, K. A. Wahid, P. Babyn, H. R. Marateb, M. Mansourian, and R. Sami, "COVID-SAFE: An IoT-based System for Automated Health Monitoring and Surveillance in Post-Pandemic Life," *IEEE Access*, vol. 8, pp. 188 538–188 551, 2020.

[17] E. Luo, M. Z. A. Bhuiyan, G. Wang, M. A. Rahman, J. Wu, and M. Atiquzzaman, "PrivacyProtector: Privacy-Protected Patient Data Collection in IoT-Based Healthcare Systems," *IEEE Communications Magazine*, vol. 56, no. 2, pp. 163–168, 2018.

[18] E. Ghasemi Komishani, M. Abadi, and F. Deldar, "PPTD: Preserving Personalized Privacy in Trajectory Data Publishing by Sensitive Attribute Generalization and Trajectory Local Suppression," *Knowledge-Based Systems*, vol. 94, pp. 43–59, 2016.

[19] L. Yao, Z. Chen, H. Hu, G. Wu, and B. Wu, "Sensitive Attribute Privacy Preservation of Trajectory Data Publishing Based on l-diversity," *Distributed and Parallel Databases*, vol. 39, no. 3, pp. 785–811, 2021. [Online]. Available: https://doi.org/10.1007/s10619-020-07318-7

[20] R. Tojiboev, W. Lee, and C. C. Lee, "Adding Noise Trajectory for Providing Privacy in Data Publishing by Vectorization," *In the Proceedings of - IEEE International Conference on Big Data and Smart Computing, BigComp 2020*, pp. 432–434, 2020.

[21] L. Yao, Y. Zhang, Z. Zheng, and G. Wu, "GAN-Based Differential Privacy Trajectory Data Publishing with Sensitive Label," *In the Proceedings of - 8th International Conference on Big Data Computing and Communications, BigCom 2022*, pp. 112–119, 2022.

[22] R. Wen, W. Cheng, H. Huang, W. Miao, and C. Wang, "Privacy Preserving Trajectory Data Publishing with Personalized Differential Privacy," *In the Proceedings of - ISPA-BDCloud-SocialCom-SustainCom 2020*, pp. 313–320, 2020.

[23] T. Kanwal, S. A. A. Shaukat, A. Anjum, S. u. R. Malik, K. K. R. Choo, A. Khan, N. Ahmad, M. Ahmad, and S. U. Khan, "Privacy-Preserving Model and Generalization Correlation Attacks for 1:M Data with Multiple Sensitive Attributes," *Information Sciences*, vol. 488, pp. 238–256, 2019. [Online]. Available: https://doi.org/10.1016/j.ins.2019.03.004

[24] T. Li, N. Li, J. Zhang, and I. Molloy, "Slicing: A New Approach for Privacy Preserving Data Publishing," *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 3, pp. 561–574, 2012.

[25] J. Jayapradha, M. Prakash, Y. Alotaibi, O. I. Khalaf, and S. A. Alghamdi, "Heap Bucketization Anonymity - An Efficient Privacy-Preserving Data Publishing Model for Multiple Sensitive Attributes," *IEEE Access*, vol. 10, pp. 28 773–28 791, 2022.

[26] L. Yao, Z. Chen, X. Wang, D. Liu, and G. Wu, "Sensitive Label Privacy Preservation with Anatomization for Data Publishing," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 2, pp. 904–917, 2021.

[27] Z. S. H. Abad, D. M. Maslove, and J. Lee, "Predicting Discharge Destination of Critically Ill Patients Using Machine Learning," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 3, pp. 827–837, 2021.

[28] F. Song, T. Ma, Y. Tian, and M. Al-Rodhaan, "A New Method of Privacy Protection: Random k-Anonymous," *IEEE Access*, vol. 7, pp. 75 434–75 445, 2019.

[29] L. Zhang, J. Xuan, R. Si, and R. Wang, "An Improved Algorithm of Individuation K-Anonymity for Multiple Sensitive Attributes," *Wireless Personal Communications*, vol. 95, no. 3, pp. 2003–2020, 2017.

[30] Z. Li and X. Ye, "Privacy Protection on Multiple Sensitive Attributes," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 4861 LNCS, pp. 141–152, 2007.

[31] V. S. Susan and T. Christopher, "Anatomisation with Slicing: A New Privacy Preservation Approach for Multiple Sensitive Attributes," *SpringerPlus*, vol. 5, no. 1, pp. 1–21, 2016.

# Machine Learning in Malware Analysis: Current Trends and Future Directions

Safa Altaha[1], Khaled Riad[2]
College of Computer Sciences & Information Technology,
King Faisal University, Al-Ahsa 31982, Saudi Arabia[1]
Computer Science Department, College of Computer Sciences & Information Technology,
King Faisal University, Al-Ahsa 31982, Saudi Arabia[2]
Mathematics Department-Faculty of Science, Zagazig University, Zagazig 44519, Egypt[2]

*Abstract*—**Malware analysis is a critical component of cyber-security due to the increasing sophistication and the widespread of malicious software. Machine learning is highly significant in malware analysis because it can process huge amounts of data, identify complex patterns, and adjust to changing threats. This paper provides a comprehensive overview of existing work related to Machine Learning (ML) methods used to analyze malware along with a description of each trend. The results of the survey demonstrate the effectiveness and importance of three trends, which are: deep learning, transfer learning, and XML techniques in the context of malware analysis. These approaches improve accuracy, interpretability, and transparency in detecting and analyzing malware. Moreover, the related challenges and Issues are presented. After identifying these challenges, we highlight future directions and potential areas that require more attention and improvement, such as distributed computing and parallelization techniques which can reduce training time and memory requirements for large datasets. Also, further investigation is needed to develop image resizing techniques to be used during the visual representation of malware to minimize information loss while maintaining consistent image sizes. These areas can contribute to the enhancement of machine learning-based malware analysis.**

*Keywords—Malware; malware analysis; machine learning; deep learning; transfer learning*

## I. Introduction

The increasing demand for computer system usage and internet connectivity results in the continuous evolution and change of malware [1]. As more people and businesses rely on computers and the internet in their daily activities, the potential attack surface for malware actors and attackers expands, giving them more opportunities to exploit vulnerabilities and compromise systems. Moreover, the widespread adoption of the internet has connected millions of devices, which allows malware to spread rapidly. Malicious actors can use the internet to distribute malware through many channels, such as: email attachments, malicious websites, or social media platforms. Malware describes any software or collection of instructions that are designed to intentionally affect computer systems, businesses, or users by causing harm. [2]. The term "malware" includes a wide range of threats, which include viruses, worms, trojan horses, ransomware, adware, and various other types of malicious software [3]. The detection and analysis of malware has consistently been a major concern and a challenging issue due to limitations in analysis methodologies, performance accuracy, and approaches that fail to identify unforeseen malware attacks. Malware analysis involves the utilization of techniques from diverse fields such as program analysis and network analysis. Its purpose is to examine malicious samples in order to gain a comprehensive understanding of various aspects, including their behavior and the transformations they undergo over time [4]. Recently, researchers have introduced various techniques for malware analysis. These techniques are typically classified into two groups: the first one is the signature-based techniques and the second one is ML-based techniques. The first techniques depend on predefined patterns or signatures of known malware samples to identify and detect malicious software. ML-based techniques for malware detection utilize ML algorithms to analyze benign and malicious malware samples. By examining these samples, the algorithms generate learning patterns that can be used to detect both known and unpredictable and new malware. This ability makes ML-based approaches a preferred choice for malware detection. ML is a core element of artificial intelligence, it typically enables systems to automatically learn and improve from experience, without requiring explicit programming for each task [5], [6]. Compared to signature-based techniques, which rely on predefined signatures, ML-based techniques are more efficient in identifying new malware [7]. This is because the accuracy of ML models depends on the features used and the training set, allowing them to learn and detect the changing characteristics of malware.

This paper aims to provide a comprehensive and up-to-date overview of existing trends related to ML used to analyze, detect, and classify malware, which includes a description of each trend, its related challenges, and issues. In addition, we include future direction suggestions. This research paper attempts to address the following questions:

1) What are the trends related to malware analysis mechanisms that use machine learning?
2) What are the potential issues and challenges related to each trend related to malware analysis mechanism?
3) What could be the future directions of research in this domain that need more investigations?

The rest of the paper is organized as follows. Section II introduces the three basic approaches to malware analysis. Section III presents the search strategy. Followed by section IV which describes different trends in malware analysis using ML. Section V states some challenges related to each trend. In

Section VI, some future directions and work are highlighted. Finally, Section VII concludes this work.

## II. MALWARE ANALYSIS APPROACHES

Three basic approaches to malware analysis and detection include static, dynamic, and hybrid analysis. All of them play significant roles in the overall malware analysis process and each contributes to the overall malware analysis process in different ways.

Static analysis involves inspecting the structure of an executable file without executing it. The executable file has various static attributes, such as distinct sections and memory compactness. The static analysis involves two parts: basic and advanced. Basic static analysis focuses on the basic properties and features of the malware to gain an initial understanding of its characteristics. File size, file type, and header information are extracted using various tools and examined during basic static analysis. After basic static analysis, advanced static analysis techniques can be applied to gain a deeper understanding of the malware's behavior and capabilities. The advanced static analysis features require more in-depth tools to extract and uncover the actual behavior of malware. Advanced static analysis involves investigating the program commends in detail [8]. The static analysis faces challenges in detecting obfuscated malware, as it is unable to effectively analyze packed samples. This limitation has been highlighted by Komatwar and Kokare [9].

Dynamic analysis involves executing the program instructions and examining the behavior of malware. To secure the machine from being affected by the malware, the dynamic analysis is conducted within an isolated environment, like a virtual machine or sandbox. Dynamic analysis can be divided into two parts: basic and advanced. In basic dynamic analysis, monitoring tools are used to examine the behaviors of malware. On the other hand, advanced dynamic analysis involves the use of debugging tools. These debuggers enable the analysts to run individual commands, with the ability to modify parameters and variables [8]. During dynamic analysis, the software operates in an environment where it has full access to all resources. At the conclusion of malware execution, the controlled environment is restored to its default state, which was captured at the beginning of the environment setup. An agent within the controlled environment logs the software's behavior [10], [11]. Unlike static analysis, dynamic analysis can deal with obfuscation and detect new malware.

Hybrid analysis is an approach that combines both static analysis and dynamic analysis techniques to detect and analyze malware. This method begins by initially analyzing the malware statically, examining its code and structure without executing it. Then, dynamic analysis is applied to enhance the overall analysis further. By incorporating dynamic analysis, the hybrid approach overcomes the limitations of applying static or dynamic analysis alone. It allows for a more comprehensive understanding of the malware's behavior. The combination of both analysis approaches enhances the overall analysis process and improves the effectiveness of malware detection and analysis [12].

Table I presents the pros and the cons of each approach [13],[14].

TABLE I. COMPARISON BETWEEN MALWARE ANALYSIS APPROACHES

| Approaches | Pros | Cons |
|---|---|---|
| Static analysis | Require less time and less power and memory consumption. | Can't detect obfuscation and unknown malware. |
| Dynamic analysis | The ability to detect unknown malware | Consumes higher amount of resources, unsafe, and require more time. |
| Hybrid analysis | Result in more accurate result. | Higher complexity and higher cost. |

Based on a study that was conducted by Gorment et al. [1] for ML algorithms for malware detection, they found out that most contributions on this domain use static analysis with a percentage of 53.3% of the related studies, while dynamic analysis accounted for 28.9% and followed by hybrid analysis with percentage 17.8%. This demonstrates the high effectiveness of static analysis as it is fast and safe. In addition, it has a small false positive rate compared to dynamic analysis [15].

## III. RESEARCH STRATEGY

This section outlines the methodology employed in this study, including the search strategy and the criteria used to determine the final set of papers included in the analysis. Google Scholar and other databases such as IEEE and ResearchGate are used to search for existing related literature. We have included research papers that are related to Malware analysis and ML. After that, the duplicated papers were excluded. Also, some papers were excluded for other reasons, such as they are not relevant, the whole paper is not available, only the abstract is provided, and some papers are written in foreign languages.

To make sure that we include recent and relevant research, we selected papers published within the past four years. Moreover, priority was given to papers that focused on ideas or objectives that have been the focus in recent years and had not been explored by researchers in previous years. The selected publications specifically emphasized contributions related to malware analysis or detection in general, with a particular focus on ML-based methods for malware analysis. By following this methodology, the goal of this study is to gather a comprehensive and up-to-date collection of literature that addresses the intersection of malware analysis and ML.

## IV. TRENDS IN MALWARE ANALYSIS USING MACHINE LEARNING

Numerous techniques, trends, and strategies have been proposed for malware analysis and detection that make use of ML. This section explores trends related to this domain that was introduced and proposed by researchers in recent years as the field of malware analysis and ML is dynamic and continuously improving.

### A. Deep Learning-based Malware Analysis

DL is an advanced subset of ML that brings ML closer to the field of artificial intelligence. It enables the modeling of complicated relationships and concepts by employing multiple layers of representation [16]. The motivation behind DL is the function and organization of the human brain and its interconnected network of neurons. It consists of multiple
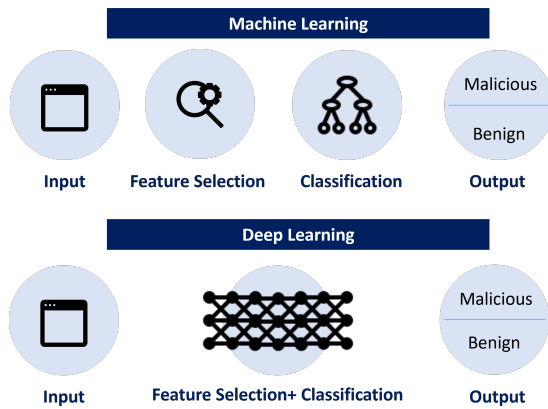
Fig. 1. Difference between deep learning and machine learning.

layers of connected nodes, called neurons or units. Each neuron takes inputs and then performs some computation to produce an output that is used as input to the next layer. The layers are organized hierarchically. In DL, each layer of a neural network learns more abstract concepts from the data. The higher layers build based on the representations learned by the lower layers. DL does not require explicit feature engineering or preprocessing of the data. Instead, it automatically extracts relevant features and representations directly from the raw data [17].

The motivation behind the adoption of DL in many fields is the need to organize and analyze massive volumes of data efficiently. Moreover, DL models have the capability to learn and extract relevant features directly from the raw input data during the training process [18], which is the main difference between ML and deep learning as shown in Fig. 1. It is commonly preferred and used in many domains, including image processing, speech processing, healthcare, and the rapidly expanding field of cybersecurity, which has seen a surge in demand for DL techniques [15].

Rhode et al. [19] proposed a Recurrent Neural Network (RNN) model for malware analysis and prediction. Their focus was on examining the possibility of predicting that an executable file is malicious or normal using a brief snapshot of behavioral data. They found out that employing an ensemble of RNN allows for accurate classification of whether an executable is malicious or benign within the initial five seconds of executing the file, achieving an impressive accuracy rate of 94%. Elayan and Mustafa [20] made use of two static features of Android applications, which are Application Program Interface (API) and permission, to present an approach for detecting malware in Android applications. Gated Recurrent Units (GRUs) were used, which is a type of RNN. GRU consisted of three blocks: input block, middle block, and output block. The approach achieved good results, correctly predicting 98.0% of the dataset. It showed high scores in both recall and precision, successfully detecting 99.2% of the malware samples.

Catak et al. [21] developed a classification method using Long Short-Term Memory (LSTM). LSTM is a DL method that was developed as an improvement over RNN. LSTM was specifically designed to address the limitations of RNN. The

malware is analyzed and detected based on the API class that the Windows operating system makes; this is used as a representation of the patterns of malicious software. The classifier's result indicates a high level of accuracy, reaching up to 95%. McDole et al. [22] introduced a study focusing on the analysis of malware detection techniques in cloud Infrastructure-as-a-Service (IaaS) environments, along with exploring the potential of using CNN. The proposed malware detection by Ravi et al. [23] used a CNN to detect malware in smart healthcare systems for both Windows and Android. The proposed approach achieved an accuracy of 98% in the Windows dataset and 97% in the Android dataset.

Bayazit et al. [24] developed comparative systems that is based on DL for malware detection, where they employed various approaches and compared the results of each approach. The system used both Static and dynamic analysis with different ML and DL classifiers. Moreover, A comparative analysis is conducted, comparing traditional ML algorithms such as Decision Trees, Random Forests, and DL algorithms including LSTM, CNN-LSTM, ANN, and Multilayer Perceptron (MLP). They found that LSTM achieved a high accuracy rate of 98.75% in static analysis classification. Additionally, the CNN-LSTM deep learning algorithm showed a high accuracy rate of 95.26% in dynamic analysis classification.

İbrahim et al. [25] introduced an approach that used static analysis and the most important features from Android applications, including two newly proposed features. These features are then used as input for an API DL model specifically developed for this purpose. The proposed method is implemented and evaluated using a classified dataset of Android applications. They focused on extracting the following features: permissions, API calls, services, broadcast receivers, and opcode sequences. Furthermore, they introduced two new static features: application size and fuzzy hash.

Patil and Deng [26] showed how accuracy could be improved using DL networks rather than the traditional ML models by introducing a neural networks-based framework for malware analysis that achieved high accuracy. The findings from the experiment indicated that the DL-based malware classification method achieves high accuracy in classification. Moreover, they suggested that the backpropagation and gradient descent mechanisms in DL help improve accuracy, true positive rate, and reduced false positive rate.

Rodrigo et al. [27] designed a hybrid ML model for Android malware detection. The model contained three fully connected neural networks: The first network focused on static features and achieved an accuracy of 92.9% when trained on 840 static attributes. The second network is designed for dynamic analysis and achieved an accuracy of 81.1% when trained on 3722 dynamic attributes. The last network combined both static and dynamic features, resulting in a hybrid model that achieved an accuracy of 91.1% when trained on 7081 static and dynamic attributes. This demonstrates that the hybrid analysis performs better than using static and dynamic features.

Obaidat et al. [28] presented an approach combining static analysis techniques with advanced image-based DL algorithms to improve the accuracy of malware detection in Java bytecode. The proposed approach, called Jadeite, is designed to classify Java bytecode files as malicious or benign using a combination

of supervised DL and static analysis. It takes as input a JAR file that contains the Java bytecode and consists of three main components. The first component is the Bytecode Transformation Engine, which converts the Java bytecode file into a grayscale image. The second component is the Feature Extraction Engine, which takes the features from the bytecode file. Finally, the CNN classifier Engine takes both the grayscale image and the features obtained from the previous steps. It employs a CNN model to detect if the file is malicious or benign.
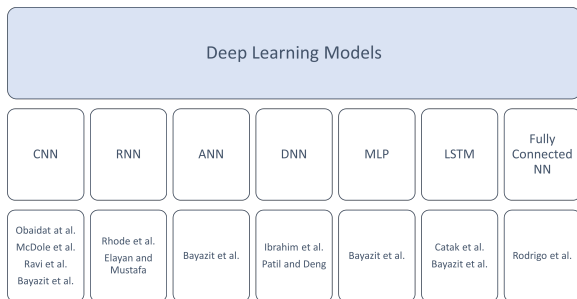


Fig. 2. Deep learning models used in the research papers.

Fig. 2 presents the deep learning techniques used by the authors in the research papers. The Figure shows that the most commonly used deep learning technique is CNN which indicates its effectiveness and efficiency in malware analysis. Following CNN, RNN is the next commonly used technique, followed by DNN and LTSM. Below is a comparative analysis these DL algorithms in malware detection:

- **CNNs** can learn the key and important features without human involvement, which is its main advantage. it is commonly used in computer vision fields for image processing and pattern recognition. Hence, it has been successful in image-based malware detection by extracting features from binary or grayscale representations of malware [29].

- **RNNs** are usually used for analyzing sequential data, making them applicable for code-based malware detection. They can capture dependencies and patterns in code instructions, enabling the identification of malicious behavior [30].

- **DNNs** also known as feedforward neural networks, can learn complex patterns and relationships in the input data, making them capable of detecting malware based on learned representations. Wang et al. [31] have stated that recent studies have shown that using DNNs in malware detection allows for the recognition of abstract patterns from extensive collections of malware samples. This enhances the ability to detect various types of malware, providing a more comprehensive approach.

- **LSTMs** are a type of RNN that can effectively capture long-term dependencies in sequential data, making them suitable for code-based malware analysis [30].

## B. Transfer Learning-based Malware Analysis

The concept of transfer learning has been explored in the literature since the 1990s, using different names such as learning to learn, life-long learning, and knowledge transfer [32]. It is a valuable technique in ML that addresses the challenge of limited training data. It allows us to use knowledge gained from a source domain and apply it to a target domain, even when the training and test data are dependent and not identically distributed. This approach is particularly helpful in domains where improving performance is difficult due to a lack of enough training data [33]. The main difference between traditional ML and transfer learning lies in the treatment of tasks and the use of previous knowledge. In traditional ML, each task is treated independently, and the model needs to learn from scratch for each task. There is no sharing or transfer of knowledge between tasks, and the model starts from the beginning for each new task. On the other hand, transfer learning involves taking knowledge and insights obtained from previous tasks or domains and transferring them to the learning process of a new task. Instead of starting from scratch, the model can benefit from the knowledge extracted from other source tasks [34]. Fig. 3 shows that the traditional ML learns the task from scratch. This means that they are trained on specific datasets related to the task. On the other hand, transfer learning uses the knowledge gained from a learning system and applies it to another learning system [35].

Transfer learning offers several advantages, including faster learning and reduced reliance on large training data. By applying knowledge from related tasks or domains, the learning process becomes quicker because the model is built based on existing knowledge. Also, the need for a large amount of training data is reduced as the model can generalize effectively by taking advantage of the existing knowledge [32].
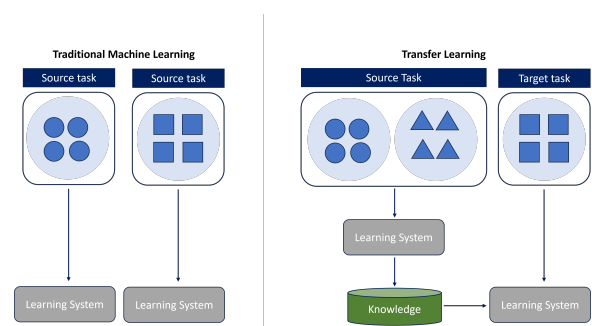


Fig. 3. Learning process of traditional ML and transfer learning.

Chen [36] introduced an approach that used deep transfer learning for static malware classification. He showed that the proposed technique has better performance compared to training from scratch and other traditional ML models. Moreover, the transfer learning scheme speeds up the training phase in deep neural networks while maintaining high performance in classification.

Bhodia et al. [37] used transfer deep learning for malware analysis and detection based on image analysis. The authors converted executable files into images and used DL models for image recognition. To train the models, they used transfer

learning by applying existing DL models that have been trained before on large image datasets. Transfer DL result was compared to a simpler k-Nearest Neighbor approach (k-NN). In certain cases, the k-NN learning technique had better performance than the proposed models. However, in simulated zero-day experiments, the proposed models had better performance compared to k-NN. Ahmad et al. [38] proposed an approach for classifying malware into nine different classes. They treated the malware binaries as image data and applied various ML and DL techniques. Logistic Regression, ANN, CNN, transfer learning on CNN, and LSTM models were used to achieve the classification results. Additionally, they employed transfer learning with InceptionV3 for training, which showed better results, particularly when compared to the LSTM model.

Acharya et al. [39] presented a transfer learning approach for efficient Android malware detection. First, the authors performed malware detection using the traditional ML models such as CNN and then they used the transfer learning technique to compare between them. They transferred the relevant features and information from a model that was trained before to a target model, and it was found that transfer learning resulted in low computational costs.

Prima and Bouhorma [40] suggested a framework for classifying malware using transfer learning. The proposed approach used pre-trained DL models that have been trained on large image datasets. In this framework, the input is a grayscale image that represents the malicious program. Then, this image is passed over a block of convolutional layers. The output is the class of the malware. According to their findings, Prima and Bouhorma stated that CNN shows better performance compared to traditional ML techniques, particularly in domains like image classification. Zhao et al. [41] proposed a malware classification method that incorporates transfer learning, multi-channel image vision features, and CNN. The methodology involves several steps. First, they extracted features from malware samples and converted them into grayscale images of three different types. To ensure a consistent size, they process the grayscale image sizes using an algorithm. Next, they synthesized the three grayscale images into three-dimensional RGB images. These RGB images are than used for training and classification.

Panda et al. [42] used two malware image datasets, namely Malevis and MalImg. For preparing the datasets, they performed pre-processing and resized each image to 224x224 pixels. Since the datasets were relatively small, they applied a technique to increase the amount of data available for training. To evaluate the accuracy of the proposed model, they used the MalImg dataset as a baseline. On the other hand, the Malevis dataset was used to build a transfer learning on CNN. The approach involved developing a transfer learning model from scratch. The extracted features are used as input to three neural network models: the autoencoder, the gated recurrent unit, and the multi-layer perceptron. The result of the multi-layer perceptron model is the final classification by taking the output from the gated recurrent unit model as its input. Tasyurek and Arslan [43] proposed a fast and accurate model that is based on CNN. The model converts the features obtained from the manifest file to an RGB format image. Then, these images are used to train the model with the transfer learning technique. The experiment results in high accuracy equal to 98.3% and quick prediction due to the use of transfer learning.

He et al. [44] claimed that the traditional ML models are not effective in accurately identifying previously unknown and zero-day malware using Hardware-Based Malware Detection. So, they proposed a Hardware-Based Malware Detection based on deep neural networks and transfer learning. the proposed solution which is based on image-based hardware events showed better performance than the existing ML methods. It achieves a high detection rate of 97% at runtime, using only the top four hardware events. Additionally, the solution maintains a low false positive rate and does not require any hardware redesign overhead. Ngo et al. [45] proposed a model that extracts both static and dynamic features for detecting malware. Then, they presented a technique for transferring knowledge that is obtained from a big source model, which is trained on the previously extracted features, to a small target model. The authors stated that the proposed model reduces the time required for prediction.

*C. Explainable Machine Learning-based Malware Analysis*

In order to create ML models that are reliable to humans, researchers have discussed various techniques for explaining and interpreting these models to users. This field of study, known as "explainability," focuses on reasoning and decision-making processes employed by ML models. Explainability includes any technique that helps users or developers understand the behavior of ML models [46].

The term "XAI" (eXplainable Artificial Intelligence) was introduced by the Defense Advanced Research Projects Agency (DARPA) in 2017. Since then, it has gained popularity across various fields, including healthcare, transportation, legal, finance, military, and scientific research. XML refers to the development and application of ML models and algorithms that are transparent, interpretable, and able to provide understandable explanations for their predictions or decisions. XML aims to bridge the gap between the complexity of traditional ML models and the need for human interpretability and understanding. In traditional ML, models such as deep neural networks can be complex and hard to understand. They operate as black boxes, making predictions without providing any explanation for the underlying reasoning or factors that affect the output. This lack of transparency can be a challenging issue, especially in domains where trust, accountability, and interpretability are important, such as healthcare, finance, and legal systems [47]. In contrast, the goal of XML techniques is to make models more transparent by providing explanations for their predictions.

Alani and Awad [48] introduced a lightweight Android malware detection system that applies of XML techniques. The system uses static features extracted from applications to classify the application as malicious or benign. The results of the experiment show an accuracy rate of over 98%, while the system remains lightweight and consumes little device's resources. Moreover, the classifier model is interpreted using Shapley Additive Explanation (SHAP) values. Liu et al. [49] also have proposed Android malware detection based on XAL. This study takes a different approach compared to other research. Instead of focusing on evaluating how well ML classifiers detect malware and identify the causes, it applies

XML approaches. The goal is to understand what ML-based models learn during training. They stated that authors should have a better understanding of how ML models work, rather than just focusing on improving the accuracy. Iadarola et al. [50] presented a DL model that is based on images for detecting which family the malware belongs to. Moreover, they explained the system prediction by using LIME as the explanation method. They achieved high accuracy equal to 93.4%.

Manthena [51] stated that many existing studies in the field of malware analysis lack transparency and explainability regarding the predictions made by their models. This absence of information about the models' decisions is a significant issue, particularly in the context of malware analysis. Manthena presented an online malware analysis by applying XML, such as KernalSHAP, TreeSHAP, and DeepSHAP to analyze and evaluate the reliability of the performance metrics. By doing so, she emphasized the importance of applying XML techniques in the context of online malware detection. While Manthena used SHAP as an XML technique, Kinkead et al. [52] used Local Interpretable Model-Agnostic Explanations (LIME) to explain the predictions of a ML model. They have developed a method that used CNN to identify significant locations within an Android app's opcode sequence. These identified locations are believed to play a key role in detecting malware. Secondly, they have conducted a comparison between the locations highlighted by CNN and the locations identified as important by LIME. This comparison allows us to evaluate the consistency between the two methods in identifying significant locations for malware detection. After conducting the experiment, they found out that the positions in the opcode sequences classified as malicious by CNN closely match those classified as most malicious by LIME.

Lu and Thing [53] proposed an android malware detection using three model explanation methods, which are Modern Portfolio Theory (MPT), SHAP, and LIME. They conducted an experiment to compare these methods regarding the explanation ability. The result of the experiment showed that the MPT is considered valuable to security analysis as it can be used to determine the reasons that the classifiers are fooled by the adversarial samples. Iadarolaa et al. [54] proposed a method for Android malware detection and family identification that depends on representing applications as images. These images are then used as input to an explainable DL model specifically designed by the authors. The purpose of this approach is to provide a transparent and interpretable solution for detecting and identifying malware in Android applications. The proposed methodology can be divided into two main parts. The first part involves training the model using appropriate techniques and data to achieve perfect performance. The second part of the methodology is responsible for the interpretability of the model's learning process. This aspect emphasizes understanding and explaining how the model makes predictions or decisions. Pan et al. [55] proposed a solution that focuses on overcoming the limitations of current malware detection methods, which include prediction inaccuracy and a lack of transparency. To address these two challenges, they have developed a hardware-assisted malware detection framework using an XML algorithm based on regression.

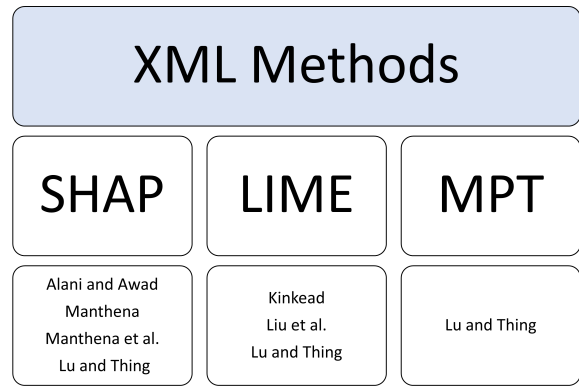Manthena et al. [56] addressed the block box characteristics



Fig. 4. XML Methods.

found in many ML and DL models such as CNN and Feed-Forward Neural Net (FFNN) by proposing a malware detection system that is trained on an online dataset and using SHAP in order to explain the outcome of the model. Sharma et al. [57] proposed an explainable system for malware detection using traffic analysis. Explainability in this model is gained by using features of the network traffic that are understandable by humans and interpretable ML decision trees for detecting malware.

Fig. 4 shows that the most commonly used methods for model explanation in malware analysis are SHAP and LIMA. The main difference between these methods is that SHAP provides a general explanation, whereas LIME produces local reasoning. This means that SHAP provides an explanation regarding the performance of the model across all samples. LIMA Interpret the prediction of a model for a specific sample[58].

Table II provides a summary of the analyzed papers for the reader and researcher to have a clear understanding.

*D. Major Findings and Discussion*

This section discusses some key findings based on the surveyed works in the above section. According to the previous section, the majority of current studies make use of DL models. The use of DL models allows for effective detection and classification of malware based on complex patterns and features. Additionally, transfer learning has been widely applied in malware analysis, using pre-trained models on large datasets to enhance the performance of malware detection systems. Moreover, XML techniques have gained attention in the field of malware analysis. XML methods provide an understanding of the decision-making process of the models, that ensure transparency and interpretability. Overall, the surveyed works show the effectiveness and importance of using DL, transfer learning, and XML techniques in the context of malware analysis.

The method of converting the samples or the files to an image before processing is used by many researchers based on surveyed work, such as [28], [37], [38], [40], [41], [42], [43], [44], [50], [54].

TABLE II. SUMMARY OF RECENT WORKS

| Authors | Pub. year | Focus/Objective | Platform | Machine-learning techniques used |
|---|---|---|---|---|
| Rhode et al. [19] | 2018 | Performing early detection of Maliciousness of a file within the initial 5 seconds of its execution. | Windows 7 Executable file | RNN |
| Elayan and Mustafa [20] | 2021 | Utilizing DL techniques, specifically the GRU architecture of Recurrent Neural Networks to detect malware in Android applications. | Android | |
| Obaidat et al. [28] | 2022 | Identifying Java bytecode malware programs through static analysis and image-based DL classification. | Java program | CNN |
| McDole et al [22] | 2021 | Analyzing malware in a cloud environment by utilizing DL. | Cloud IaaS | |
| Ravi et al. [23] | 2022 | Proposing a method for malware detection on Windows and Android operating system in smart healthcare systems. | Windows and Android OS | |
| Catak et al. [21] | 2020 | Using LSTM for malware analysis based on API calls in Windows operating systems. | Windows OS | LSTM |
| Bayazit et al. [24] | 2023 | Present and compare the use of static and dynamic analysis in DL-based malware binary classification. | Android | LSTM, CNN-LSTM, ANN and MLP . |
| Rodrigo et al. [27] | 2021 | Using hybrid ML for Android malware detection. | Android | Fully connected neural networks |
| İbrahim et al. [25] | 2022 | Using static analysis for android malware detection. | Android | Deep learning |
| Patil and Deng [26] | 2020 | Comparing the accuracy of the DL approach and traditional ML in malware analysis. | - | |
| Kinkead et al. [52] | 2021 | Using Explainable CNNs for developing Android Malware Detection method. | Android | . |
| ladarola et al. [50] | 2023 | Explinable DL model on images for Malware fimaly detection. | - | Explainable CNNs |
| Pan et al. [55] | 2020 | Addressing two limitations: inaccuracy in predictions and lack of transparency | - | Explainable RNN |
| Liu et al. [49] | 2022 | Researchers should focus on having a better understanding of how ML models work, rather than just focusing on improving the accuracy of the models. | Android | Explainable machine learning |
| Manthena [51] | 2022 | Addressing the lack of Explainable ML approaches for online malware analysis. | - | |

| | | | | |
|---|---|---|---|---|
| Alani and Awad [48] | 2022 | Explaining the reasons behind the selected features using Shapley additive explanation values to ensure that the high accuracy of the classifier come from explainable conditions. | Android | Explainable machine learning |
| Manthena et al. [56] | 2023 | Addressing the block box issue that found in many Ml and DL models | - | |
| Lu and Thing [53] | 2022 | Designing Android malware analysis and Comparing the explanation between different explanation methods. | Android | |
| Sharma et al. [57] | 2023 | Explainable malware detection using network traffic features | - | Explainable Decision tree |
| Iadarolaa et al. [54] | 2021 | Using explainable deep learning model for detecting malware and identifying their family. | mobile | Explainable deep learning |
| Chen [36] | 2019 | Demonstrating that transfer learning outperforms training from scratch. | - | Transfer deep learning |
| Bhodia et al. [37] | 2019 | Malware detection and classification based on image analysis | Executable files | |
| Prima and Bouhorma [40] | 2020 | CNN have better performance compared to traditional learning techniques. | - | |
| Panda et al. [42] | 2023 | Developing a transfer learning model from scratch | IoT | Transfer learning on CNN |
| Acharya et al. [39] | 2023 | Using transfer learning for mobile malware detection to lower the computational cost. | Mobile | |
| Zhao et al. [41] | 2023 | Solving the problem of existing DL model for malware detection which is long training time | - | |
| Ahmad et al. [38] | 2023 | Classifying malware, with nine different class using CNN-transfer learning model | - | |
| Ngo et al. [45] | 2023 | Addressing the limitations of feature aggregation while using static and dynamic features and transferring knowledge from big to small models. | - | |
| Tasyurek and Arslan [43] | 2023 | Fast CNN-based transfer learning for malware analysis | Android | |
| He et al. [44] | 2022 | Using a deep neural network and transfer learning for detecting zero-day malware. | - | Transfer learning on DNN |

Usually, the malware sample comes as an executable file in binary format. To perform malware representation, the byte in the binary file is converted to pixels in the image containing textural patterns which results in better visualization of the malware [59], [60]. Malware representation of a file before processing it is very common in malware analysis files because by using this approach, researchers avoid the need for and dependency on features engineering process and methods [61], [62]. Moreover, it allows for a visual representation of the malware and incorporates computer vision with this domain [63].

## V. LIMITATIONS AND CHALLENGES

Over time, various trends in malware detection and analysis have emerged, each with its limitations and drawbacks. As a result, researchers have shifted their focus to alternative technologies that aim to detect malware in real-time with minimal false positives and increase detection and classification accuracy. This section discusses the challenges and limitations associated with each trend based on the analyzed papers in the previous sections.

### A. Limitations and Challenges Found in Papers using Deep Learning

- Rhode et al. [19] used RNN for early-stage malware detection assuming that the malware will execute the malicious activities within the first five seconds. But what if the attacker or adversaries are aware that the first 5 seconds of a file's execution are used to determine whether it is malicious or not, they can manipulate the file's behavior to avoid detection. By introducing long periods of sleep or inactivity at the beginning of a malicious file, so, the attacker can trick the system into classifying the file as safe. We suggest incorporating additional features or mechanisms in the system to capture and analyze behavior beyond the initial 5 seconds, such as by extracting and analyzing a broader range of features, like system calls, or network activity to gain a deeper understanding of the file's intent.

- Obaidat et al. [28] used DL with a large dataset for Java malware detection. DL architectures rely heavily on supervised learning, which requires a large amount of labeled samples to train the model effectively. As mentioned by Kadam and Vaidya in [64], using a small number of samples does not allow the model to learn the underlying features accurately in the training stage, and processing such a large amount of data in DL requires extensive training time and needs high processing power.

- The dataset used by Catak et al. [21] has an unequal distribution of instances for each malware Category. For example, there are 1001 rows labeled as Worms and only 379 samples for Adware. This may affect the performance of the model in classification [65]. To address this issue, we can use data resampling techniques. Data resampling is used for class-imbalance problems which aim to balance the distribution of instances across different classes by either increasing the number of instances in the minority class (over-sampling) or decreasing the number of instances in the majority class (undersampling). This helps to mitigate the impact of class imbalance on the model's training process [66].

### B. Limitations and Challenges Found in Papers using Transfer Learning

- In the work of Panda et al. [42], the transfer learning model may struggle to handle malware images with varying sizes. Inconsistencies in image sizes can lead to challenges in the model's ability to learn meaningful features and patterns. Resizing or standardizing the images to a fixed size is typically required as a pre-processing step, which can add complexity and potential loss of information. Moreover, the model proposed in this paper is incapable of detecting unseen malware. this may lead to a prolonged classification process because of the need for pre-processing the malware images before classifying the malware.

- Niu et al. [67] and Pan and Yang [68] mentioned that transfer learning suffers from an issue called negative transfer. Negative transfer refers to a phenomenon in which the application of transfer learning leads to a decrease in performance or some impact on the target task. It occurs when the knowledge learned from the source domain is not relevant or compatible with the target task. So, instead of benefiting the target task, the transferred knowledge may introduce noise or incorrect assumptions that affect the model's capability to generalize and make accurate predictions. We can mitigate the risk of negative transfer by using an ensemble of models trained with different source domains. By combining different sources of transferred knowledge, the ensemble can take advantage of the strengths of each model and reduce the impact of negative transfer.

- Based on paper of Prima and Bouhorma [40] and Panda et al. [42], We can conclude that transfer learning requires large datasets. Even though the use of a large dataset can achieve high accuracy, it is time-consuming and may require more computational resources [69]. Incremental learning could partially solve this problem. Instead of training the model on the entire large dataset at once, incremental learning can be employed to train the model on smaller subsets of data sequentially, gradually increasing the complexity of the task. This way, computational resources can be used more efficiently, and the time required for training can be reduced.

### C. Limitations and Challenges Found in Papers using XML

- The work proposed by Alani and Awad [48] has been acknowledged for its accuracy and efficiency, requiring minimal time for malware detection. However, a limitation of this approach is its inability to effectively detect unknown and obfuscated malware. This limitation arises from the fact that the proposed method relies on static features for malware detection. Static features typically refer to characteristics

extracted from the file or code without considering dynamic behaviors or run-time information. Applying dynamic analysis along with static analysis can provide more valuable information regarding the behavior of malware during execution because dynamic analysis can capture actions such as file system modifications, network communications, process interactions, and system calls, which can help identify malicious activities and detect previously unknown or obfuscated malware.

- In many cases, there exists a trade-off between model performance and interpretability as stated by Antoniadi et al. [70] and Arrietaet al.[71]. This means that as models become more interpretable, their predictive accuracy may decrease. Conversely, highly accurate models may be less interpretable. making a balance between interpretability and performance is a challenge in XML, as it requires finding the right level of transparency without decreasing the accuracy.

- The split of the dataset in Kinkead et al. [52] for training, validation, and testing was not fair. A balanced split is typically recommended to make sure that both the training and testing samples have a representative distribution of malware and benign samples. the best resampling technique for this case is cross-validation which is considered the standard resampling technique for splitting the dataset into training, validation, and testing sets [72].

- Unlike accuracy or precision, interpretability lacks standard evaluation metrics. Measuring the quality of explanations is still an active area of research. The absence of widely accepted evaluation criteria makes it challenging to compare and evaluate different XML methods [73] [74].

## VI. Future Directions

In the previous section, several issues and limitations that need to be addressed to develop a system that has the ability to detect, analyze, and classify malware efficiently are listed. In this section, we will highlight some areas for future work that researchers can use to help mitigate the mentioned issues:

### A. *The use of a Moderate-sized and Balanced Dataset*

As mentioned in the previous section some research papers used very large datasets and other papers suffered from unbalanced datasets. This problem is because the researchers believe that larger datasets lead to higher accuracy and reduce bias. While this belief is common, it is important to consider that there are potential challenges linked with very large datasets. By working with a dataset of moderate size, computational resources and time required for training models can be reduced. Another issue that needs to be taken care of is the use of up-to-date datasets because it is important for maintaining the relevance and applicability of the study in real-time scenarios.

### B. *Domain Distance Measure*

Addressing the issue of negative transfer in transfer learning is an important research challenge that needs to be

addressed in the future. Negative transfer leads to worse performance compared to not transferring at all. Therefore, it is essential to find ways to prevent negative transfer from happening in transfer learning. An accurate measurement of domain distance is another important research aspect that needs attention to solve this problem. When applying transfer learning techniques, it is essential to rate the similarity or dissimilarity between the source and target domains correctly. Existing approaches for measuring domain distance often rely on assumptions that may not capture the true underlying relationships between domains. This can cause inappropriate transfer of knowledge and performance degrading. Developing robust and accurate methods for evaluating domain distance will enable better identification of relevant knowledge for transfer and make the adaptation of models to new domains more effective [67].

### C. *Image Resizing and Standardization*

Resizing images to a fixed size is a common preprocessing step to address inconsistencies in image sizes. This step ensures that the input images have a consistent size and format, which is essential for many ML models. However, it's important to carefully consider the impact of resizing on the loss of information [75]. So further investigation regarding resizing techniques to minimize information loss while maintaining a consistent and reasonable image size for efficient processing is needed.

### D. *Distributed Computing and Parallelization*

To facilitate the training process and reduce the time required for large dataset processing, researchers should consider making use of distributed computing techniques and parallelization. This involves distributing the workload across multiple computing resources, such as GPUs or multiple machines, to train the model in parallel. Parallelization can significantly reduce the training time and reduce the memory resource requirements [76]

### E. *Comparative Evaluations*

We recommend more future research regarding evaluating the quality of the explanation of an XML. One way to solve this is by comparative evaluations. It means comparing different XML methods on fixed benchmarks or datasets. By applying multiple XML methods to the same dataset or task, their performance in terms of interpretability can be compared. This can be done using qualitative analyses by experts or by developing quantitative metrics that consider different elements of interpretability.

## VII. Conclusions

Malware analysis is a critical component of cybersecurity due to the increasing sophistication and the widespread of malicious software. Understanding malware is key to developing strong defenses. Malware analysis helps identify and classify different types of malware, which makes it easier to detect and prevent future attacks. ML plays a key role in malware analysis due to its ability to analyze large amounts of data and detect complex patterns. In this paper, we provide a survey of existing trends related to malware analysis using

ML including a description of each trend. The surveyed works show the effectiveness and importance of applying DL, transfer learning, and XML techniques in the context of malware analysis. These approaches contribute to improved accuracy, interpretability, and transparency in detecting and analyzing malware. Moreover, the challenges and limitations related to each trend are explored. Based on the survey results, we also provide some future directions to be investigated that have the potential to shape the future of malware analysis. These areas may offer exciting opportunities for further improvement in the field in order to overcome the challenges faced by the researchers.

Despite the valuable insights this paper provides, it is important to acknowledge its limitations. The sample size was relatively small, which may limit the generalizability of the findings. Future research should aim to replicate these findings with larger samples.

## REFERENCES

[1] N. Z. Gorment, A. Selamat, L. K. Cheng, and O. Krejcar, "Machine learning algorithm for malware detection: Taxonomy, current challenges and future directions," *IEEE Access*, 2023.

[2] U. V. Nikam and V. M. Deshmuh, "Performance evaluation of machine learning classifiers in malware detection," in *2022 IEEE International Conference on Distributed Computing and Electrical Circuits and Electronics (ICDCECE)*. IEEE, 2022, pp. 1–5.

[3] M. S. Akhtar and T. Feng, "Malware analysis and detection using machine learning algorithms," *Symmetry*, vol. 14, no. 11, p. 2304, 2022.

[4] D. Ucci, L. Aniello, and R. Baldoni, "Survey of machine learning techniques for malware analysis," *Computers & Security*, vol. 81, pp. 123–147, 2019.

[5] I. H. Sarker, A. Kayes, S. Badsha, H. Alqahtani, P. Watters, and A. Ng, "Cybersecurity data science: an overview from machine learning perspective," *Journal of Big data*, vol. 7, pp. 1–29, 2020.

[6] I. Sarker, "Machine learning: algorithms, real-world applications and research directions. sn comput sci 2: 160," 2021.

[7] M. T. Ahvanooey, Q. Li, M. Rabbani, and A. R. Rajput, "A survey on smartphones security: software vulnerabilities, malware, and attacks," *arXiv preprint arXiv:2001.09406*, 2020.

[8] Ö. Aslan and A. A. Yilmaz, "A new malware classification framework based on deep learning algorithms," *Ieee Access*, vol. 9, pp. 87 936–87 951, 2021.

[9] R. Komatwar and M. Kokare, "Retracted article: a survey on malware detection and classification," *Journal of Applied Security Research*, vol. 16, no. 3, pp. 390–420, 2021.

[10] M. Ijaz, M. H. Durad, and M. Ismail, "Static and dynamic malware analysis using machine learning," in *2019 16th International bhurban conference on applied sciences and technology (IBCAST)*. IEEE, 2019, pp. 687–691.

[11] J. Lee, H. Jang, S. Ha, and Y. Yoon, "Android malware detection using machine learning with feature selection based on the genetic algorithm," *Mathematics*, vol. 9, no. 21, p. 2813, 2021.

[12] N. Tarar, S. Sharma, and C. R. Krishna, "Analysis and classification of android malware using machine learning algorithms," in *2018 3rd International Conference on Inventive Computation Technologies (ICICT)*. IEEE, 2018, pp. 738–743.

[13] N. K. Gyamfi and E. Owusu, "Survey of mobile malware analysis, detection techniques and tool," 11 2018, pp. 1101–1107.

[14] V. Rao and K. Hande, "A comparative study of static, dynamic and hybrid analysis techniques for android malware detection," *International Journal of Engineering Development and Research*, vol. 5, no. 2, pp. 1433–1436, 2017.

[15] U.-e.-H. Tayyab, F. B. Khan, M. H. Durad, A. Khan, and Y. S. Lee, "A survey of the recent trends in deep learning based malware detection," *Journal of Cybersecurity and Privacy*, vol. 2, no. 4, pp. 800–829, 2022.

[16] N. Shone, T. N. Ngoc, V. D. Phai, and Q. Shi, "A deep learning approach to network intrusion detection," *IEEE transactions on emerging topics in computational intelligence*, vol. 2, no. 1, pp. 41–50, 2018.

[17] N. Rusk, "Deep learning," *Nature Methods*, vol. 13, no. 1, pp. 35–35, 2016.

[18] Karthick Jonagadla, "Deep learning in financial markets," https://www.quantace.in/deep-learning-application-financial-markets/, 2023, accessed: November 9, 2023.

[19] M. Rhode, P. Burnap, and K. Jones, "Early-stage malware prediction using recurrent neural networks," *computers & security*, vol. 77, pp. 578–594, 2018.

[20] O. N. Elayan and A. M. Mustafa, "Android malware detection using deep learning," *Procedia Computer Science*, vol. 184, pp. 847–852, 2021.

[21] F. O. Catak, A. F. Yazı, O. Elezaj, and J. Ahmed, "Deep learning based sequential model for malware analysis using windows exe api calls," *PeerJ Computer Science*, vol. 6, p. e285, 2020.

[22] A. McDole, M. Gupta, M. Abdelsalam, S. Mittal, and M. Alazab, "Deep learning techniques for behavioral malware analysis in cloud iaas," *Malware analysis using artificial intelligence and deep learning*, pp. 269–285, 2021.

[23] V. Ravi, M. Alazab, S. Selvaganapathy, and R. Chaganti, "A multi-view attention-based deep learning framework for malware detection in smart healthcare systems," *Computer Communications*, vol. 195, pp. 73–81, 2022.

[24] E. Calik Bayazit, O. Koray Sahingoz, and B. Dogan, "Deep learning based malware detection for android systems: A comparative analysis," *Tehnički vjesnik*, vol. 30, no. 3, pp. 787–796, 2023.

[25] M. İbrahim, B. Issa, and M. B. Jasser, "A method for automatic android malware detection based on static analysis and deep learning," *IEEE Access*, vol. 10, pp. 117 334–117 352, 2022.

[26] R. Patil and W. Deng, "Malware analysis using machine learning and deep learning techniques," in *2020 SoutheastCon*, vol. 2. IEEE, 2020, pp. 1–7.

[27] C. Rodrigo, S. Pierre, R. Beaubrun, and F. B. El Khoury, "A hybrid machine learning-based malware detection model for android devices. electronics 2021, 10, 2948," 2021.

[28] I. Obaidat, M. Sridhar, K. M. Pham, and P. H. Phung, "Jadeite: A novel image-behavior-based approach for java malware detection using deep learning," *Computers & Security*, vol. 113, p. 102547, 2022.

[29] L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, Y. Duan, O. Al-Shamma, J. Santamaría, M. A. Fadhel, M. Al-Amidie, and L. Farhan, "Review of deep learning: Concepts, cnn architectures, challenges, applications, future directions," *Journal of big Data*, vol. 8, pp. 1–74, 2021.

[30] P. Maniriho, A. N. Mahmood, and M. J. M. Chowdhury, "A survey of recent advances in deep learning models for detecting malware in desktop and mobile platforms," *arXiv preprint arXiv:2209.03622*, 2022.

[31] Q. Wang, W. Guo, K. Zhang, A. G. Ororbia, X. Xing, X. Liu, and C. L. Giles, "Adversary resistant deep neural networks with an application to malware detection," in *Proceedings of the 23rd ACM sigkdd international conference on knowledge discovery and data mining*, 2017, pp. 1145–1153.

[32] R. Ribani and M. Marengoni, "A survey of transfer learning for convolutional neural networks," in *2019 32nd SIBGRAPI conference on graphics, patterns and images tutorials (SIBGRAPI-T)*. IEEE, 2019, pp. 47–57.

[33] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu, "A survey on deep transfer learning," in *Artificial Neural Networks and Machine Learning–ICANN 2018: 27th International Conference on Artificial Neural Networks, Rhodes, Greece, October 4-7, 2018, Proceedings, Part III 27*. Springer, 2018, pp. 270–279.

[34] H. M. K. Barznji, "Transfer learning as new field in machine learning," 2020.

[35] M. Ranaweera and Q. H. Mahmoud, "Virtual to real-world transfer learning: A systematic review," *Electronics*, vol. 10, no. 12, 2021. [Online]. Available: https://www.mdpi.com/2079-9292/10/12/1491

[36] L. Chen, "Deep transfer learning for static malware classification," *arXiv preprint arXiv:1812.07606*, 2018.

[37] N. Bhodia, P. Prajapati, F. Di Troia, and M. Stamp, "Transfer learning for image-based malware classification," *arXiv preprint arXiv:1903.11551*, 2019.

[38] M. Ahmed, N. Afreen, M. Ahmed, M. Sameer, and J. Ahamed, "An inception v3 approach for malware classification using machine learning and transfer learning," *International Journal of Intelligent Networks*, vol. 4, pp. 11–18, 2023. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2666603022000252

[39] S. Acharya, U. Rawat, and R. Bhatnagar, "A computationally inexpensive method based on transfer learning for mobile malware detection," in *Proceedings of Fourth International Conference on Computer and Communication Technologies: IC3T 2022*. Springer, 2023, pp. 263–274.

[40] B. Prima and M. Bouhorma, "Using transfer learning for malware classification," *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 44, pp. 343–349, 2020.

[41] Z. Zhao, S. Yang, and D. Zhao, "A new framework for visual classification of multi-channel malware based on transfer learning," *Applied Sciences*, vol. 13, no. 4, 2023. [Online]. Available: https://www.mdpi.com/2076-3417/13/4/2484

[42] P. Panda, O. K. CU, S. Marappan, S. Ma, and D. Veesani Nandi, "Transfer learning for image-based malware detection for iot," *Sensors*, vol. 23, no. 6, p. 3253, 2023.

[43] M. Tasyurek and R. S. Arslan, "Rt-droid: a novel approach for real-time android application analysis with transfer learning-based cnn models," *Journal of Real-Time Image Processing*, vol. 20, no. 3, pp. 1–17, 2023.

[44] Z. He, A. Rezaei, H. Homayoun, and H. Sayadi, "Deep neural network and transfer learning for accurate hardware-based zero-day malware detection," in *Proceedings of the Great Lakes Symposium on VLSI 2022*, 2022, pp. 27–32.

[45] M. V. Ngo, T. Truong-Huu, D. Rabadi, J. Y. Loo, and S. G. Teo, "Fast and efficient malware detection with joint static and dynamic features through transfer learning," in *International Conference on Applied Cryptography and Network Security*. Springer, 2023, pp. 503–531.

[46] U. Bhatt, A. Xiang, S. Sharma, A. Weller, A. Taly, Y. Jia, J. Ghosh, R. Puri, J. M. Moura, and P. Eckersley, "Explainable machine learning in deployment," in *Proceedings of the 2020 conference on fairness, accountability, and transparency*, 2020, pp. 648–657.

[47] X. Zhong, B. Gallagher, S. Liu, B. Kailkhura, A. Hiszpanski, and T. Y.-J. Han, "Explainable machine learning in materials science," *npj Computational Materials*, vol. 8, no. 1, p. 204, 2022.

[48] M. M. Alani and A. I. Awad, "Paired: An explainable lightweight android malware detection system," *IEEE Access*, vol. 10, pp. 73 214–73 228, 2022.

[49] Y. Liu, C. Tantithamthavorn, L. Li, and Y. Liu, "Explainable ai for android malware detection: Towards understanding why the models perform so well?" in *2022 IEEE 33rd International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 2022, pp. 169–180.

[50] G. ladarola, F. Mercaldo1, F. Martinelli, and A. Santone, "Assessing deep learning predictions in image-based malware detection with activation maps," in *Security and Trust Management: 18th International Workshop, STM 2022, Copenhagen, Denmark, September 29, 2022, Proceedings*, vol. 13867. Springer Nature, 2023, p. 104.

[51] H. Manthena, "Explainable machine learning based malware analysis," Ph.D. dissertation, North Carolina Agricultural and Technical State University, 2022.

[52] M. Kinkead, S. Millar, N. McLaughlin, and P. O'Kane, "Towards explainable cnns for android malware detection," *Procedia Computer Science*, vol. 184, pp. 959–965, 2021.

[53] Z. Lu and V. L. Thing, ""how does it detect a malicious app?" explaining the predictions of ai-based malware detector," in *2022 IEEE 8th Intl Conference on Big Data Security on Cloud (BigDataSecurity), IEEE Intl Conference on High Performance and Smart Computing,(HPSC) and IEEE Intl Conference on Intelligent Data and Security (IDS)*. IEEE, 2022, pp. 194–199.

[54] G. Iadarola, F. Martinelli, F. Mercaldo, and A. Santone, "Towards an interpretable deep learning model for mobile malware detection and family identification," *Computers & Security*, vol. 105, p. 102198, 2021.

[55] Z. Pan, J. Sheldon, and P. Mishra, "Hardware-assisted malware detection using explainable machine learning," in *2020 IEEE 38th International Conference on Computer Design (ICCD)*. IEEE, 2020, pp. 663–666.

[56] H. Manthena, J. C. Kimmel, M. Abdelsalam, and M. Gupta, "Analyzing and explaining black-box models for online malware detection," *IEEE Access*, vol. 11, pp. 25 237–25 252, 2023.

[57] Y. Sharma, S. Birnbach, and I. Martinovic, "Radar: a ttp-based extensible, explainable, and effective system for network traffic analysis and malware detection," 2023.

[58] K. Safjan, "Explaining ai - the key differences between lime and shap methods," *Krystian's Safjan Blog*, 2023.

[59] "An enhancement for image-based malware classification using machine learning with low dimension normalized input images," *Journal of Information Security and Applications*, vol. 69, p. 103308, 2022.

[60] A. Bensaoud, N. Abudawaood, and J. Kalita, "Classifying malware images with convolutional neural network models," *International Journal of Network Security*, vol. 22, no. 6, pp. 1022–1031, 2020.

[61] M. U. Demirezen, "Image based malware classification with multimodal deep learning," *International Journal of Information Security Science*, vol. 10, no. 2, pp. 42–59, 2021.

[62] Z. Zhao, D. Zhao, S. Yang, L. Xu *et al.*, "Image-based malware classification method with the alexnet convolutional neural network model," *Security and Communication Networks*, vol. 2023, 2023.

[63] B. Saridou, I. Moulas, S. Shiaeles, and B. Papadopoulos, "Image-based malware detection using & alpha;-cuts and binary visualisation," *Applied Sciences*, vol. 13, no. 7, 2023. [Online]. Available: https://www.mdpi.com/2076-3417/13/7/4624

[64] S. Kadam and V. Vaidya, "Review and analysis of zero, one and few shot learning approaches," in *Intelligent Systems Design and Applications: 18th International Conference on Intelligent Systems Design and Applications (ISDA 2018) held in Vellore, India, December 6-8, 2018, Volume 1*. Springer, 2020, pp. 100–112.

[65] V. Khullar, M. Angurala, K. D. Singh, P. Prasant, V. Pabbi, and M. Veeramanickam, "Exploring methods for dealing with class imbalances in supervised machine learning structured datasets," in *2023 3rd International Conference on Advances in Computing, Communication, Embedded and Secure Systems (ACCESS)*. IEEE, 2023, pp. 209–214.

[66] A. Estabrooks, T. Jo, and N. Japkowicz, "A multiple resampling method for learning from imbalanced data sets," *Computational intelligence*, vol. 20, no. 1, pp. 18–36, 2004.

[67] S. Niu, Y. Liu, J. Wang, and H. Song, "A decade survey of transfer learning (2010–2020)," *IEEE Transactions on Artificial Intelligence*, vol. 1, no. 2, pp. 151–166, 2020.

[68] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.

[69] H. O. Ikromovich and B. B. Mamatkulovich, "Facial recognition using transfer learning in the deep cnn," *Open Access Repository*, vol. 4, no. 3, pp. 502–507, 2023.

[70] A. M. Antoniadi, Y. Du, Y. Guendouz, L. Wei, C. Mazo, B. A. Becker, and C. Mooney, "Current challenges and future opportunities for xai in machine learning-based clinical decision support systems: a systematic review," *Applied Sciences*, vol. 11, no. 11, p. 5088, 2021.

[71] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins *et al.*, "Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai," *Information fusion*, vol. 58, pp. 82–115, 2020.

[72] B. Vrigazova, "The proportion for splitting data into training and test set for the bootstrap in classification problems," *Business Systems Research: International Journal of the Society for Advancing Innovation and Research in Economy*, vol. 12, no. 1, pp. 228–242, 2021.

[73] A. Das and P. Rad, "Opportunities and challenges in explainable artificial intelligence (xai): A survey," *arXiv preprint arXiv:2006.11371*, 2020.

[74] L. Longo, R. Goebel, F. Lecue, P. Kieseberg, and A. Holzinger, "Explainable artificial intelligence: Concepts, applications, research challenges and visions," in *International cross-domain conference for machine learning and knowledge extraction*. Springer, 2020, pp. 1–16.

[75] J. J. Luke, R. Joseph, and M. Balaji, "Impact of image size on accuracy and generalization of convolutional neural networks," *Int. J. Res. Anal. Rev.(IJRAR)*, vol. 6, no. 1, pp. 70–80, 2019.

[76] K. Zhu, H. Wang, H. Bai, J. Li, Z. Qiu, H. Cui, and E. Chang, "Parallelizing support vector machines on distributed computers," *Advances in neural information processing systems*, vol. 20, 2007.

# Towards a Continuous Temporal Improvement Approach for Real-Time Business Processes

Asma Ouarhim[1], Karim Baïna[2], Brahim Elbhiri[3]
Alqualsadi research team Rabat IT Center ENSIAS Mohammed V University Rabat, Morocco[1,2]
SMARTiLAB Laboratory, Moroccan school of Engineering Sciences (EMSI) Rabat, Morocco[1,3]

*Abstract*—**Time is relative, which makes the interaction so sensitive. Indeed, contemplating the concept of real-time enterprises resembled envisioning an idealized notion that seemed unattainable and impracticable in reality. Consequently, we give a new definition of the real-time concept according to our needs and targets for a successful business process. According to this definition, we can go towards a real-time business process validation algorithm, which has the goal of ensuring quality in terms of time, i.e., time latency $\simeq 0$. Put simply, it serves as a method to assess the consistency of a process. This approach aids in comprehending the temporal patterns inherent in a process as it evolves, empowering decision-makers to glean insights and swiftly form initial judgments for effective problem-solving and the identification of appropriate solutions. Thus, our main purpose is to deliver the right information and knowledge to the right person at the right time. To achieve this, we introduce a novel real-time component within the Business Process Management Notation (BPMN), encompassing various attributes that facilitate process monitoring. This extension transforms the BPMN into a unified real-time business process meta-model. To be more specific, our contribution proposes a continuous temporal improvement assessment and knowledge management as temporal knowledge helps to evaluate the real-time situation of the business process.**

*Keywords*—*Real-time business process; real-time enterprises; temporal latency; process validation; continuous improvement approach*



Fig. 1. Concept position: real-time enterprise/process, near-real-time entreprise/process, right-real-time enterprise/process.

## I. INTRODUCTION

Business process management is one of the top development priorities in organizations; therefore, improving it becomes a priority, especially through the continuous improvement capability process [1], [2], [3]. Enhancing business process management is crucial for organizations to optimize their operations and achieve higher efficiency[1]. Our interest is time sensitivity in processes, or real-time processes, which we call right-real-time (as we will see in the research approach). Real-time enterprises entail immediate responsiveness to business demands, but in practice, achieving such instantaneous reactivity is not feasible; we are rather 'near real-time'; consequently, we depend here on customer needs that we're trying to meet through services. If an event that happens an hour from now is judged acceptable, that occurrence is now practically the standard for what constitutes real-time, in other words, right-real-time, which generates automatically time latency. One of the significant bases of our study is time latency to eliminate waste of time and have control over the whole process (see Fig. 1). Our approach aims to introduce a novel measure of capability specifically related to time. However, it is important to distinguish between two types of capabili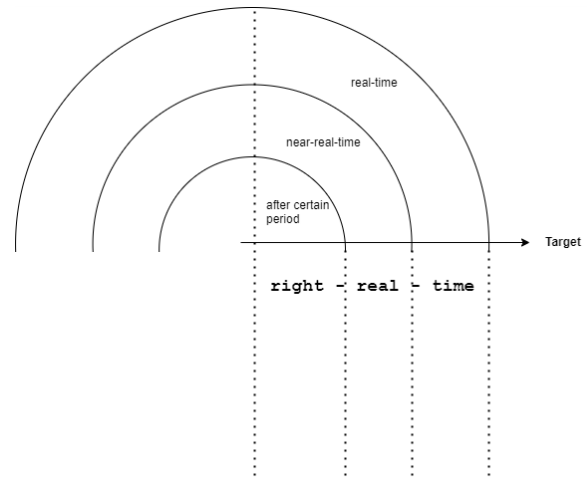ties: those that enhance an organization's ability to run processes and those that pertain to conducting business process management (BPM) [1]. Our contribution falls within the latter category. Successful Business Process Management initiatives rely on various capability factors that significantly impact their outcomes. Our objective is to define a new capability factor, namely the temporal capability factor, which plays a crucial role in the continuous improvement process. By focusing on prevention rather than cure, particularly when dealing with sensitive parameters, we can implement appropriate solutions to proactively control the situation. The continuous improvement process enables the ongoing refinement of processes and the optimization of working conditions, ultimately leading to waste elimination. (Please see the general process of continuous improvement steps below in Fig. 2, as it is inspired from [4]). After conducting extensive studies, we have discovered that time wastage has emerged as a significant concern in today's highly competitive landscape, but there is no direct tool or method that can show the real-time situation of a business process with a continuous temporal improvement method. This explains the originality of our approach, which is useful for every business process because it includes the time aspect, which allows them to identify and rectify any potential deviations, bottlenecks, or errors as they occur, preventing negative impacts on overall operations in terms of time. Effectively managing and controlling time has become a formidable challenge in the current business environment. By continually improving their Business Process Management temporal practices, businesses can streamline workflows, reduce bottlenecks,
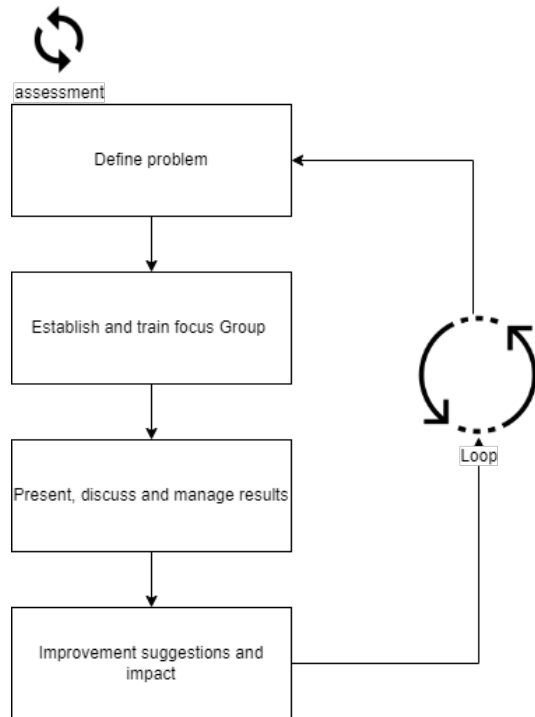
Fig. 2. CI Process.

and enhance productivity. This involves analyzing existing processes, identifying areas for improvement, and implementing strategies to enhance process efficiency and effectiveness. By focusing on Business Process Management improvement, organizations can achieve better resource allocation, reduced costs, improved customer satisfaction, and increased overall competitiveness in the market. Ultimately, continuous business process management temporal improvement leads to enhanced agility and adaptability, enabling organizations to navigate the ever-evolving business landscape with confidence [1], [5], [2], [3]. In this paper, we present a new "real-time process continuous improvement methodology" plus a "real-time process validation algorithm", which is an original search in terms of definition, modeling, and application. "Real-time Process Continuous Improvement" is an important topic for improving organizations' systems. This allows them to identify and rectify any potential deviations, bottlenecks, or errors as they occur, preventing negative impacts on overall operations in terms of time. Accordingly, we divide our contributions into four sections: Section I focuses on an overview of related works, which provides an in-depth exploration of relevant works in the field; Section II presents our proposed approach, elucidating its key components and methodologies; and Section III offers architectural thinking, which presents our approach within the enterprise architecture. Lastly, Section IV unfolds an in-depth examination, presenting the analysis, results, and discussions within the context of a case study format.

Foreword: Considering the need to establish clarity and avoid any potential confusion with existing definitions of real-time, particularly within the context of enterprise management case studies, we introduce a novel term in alignment with our specific understanding: right-real-time.

## II. RELATED WORKS

The existing literature in this field can be classified into three main categories: the advancement of business process management, the exploration of real-time enterprises, and the multifaceted understanding of time and real-time figures. These categories encompass a wide range of research and practical applications, each shedding light on different aspects of achieving efficiency and agility in organizational operations. By categorizing the related works, we can gain a comprehensive understanding of the diverse perspectives and approaches taken in the study of this subject matter.

### A. Business Process Management (BPM)

The development of Business Process Management (BPM) brings forth numerous benefits for organizations. Firstly, it enables companies to enhance their operational efficiency by streamlining and optimizing their processes, effectively eliminating bottlenecks and unnecessary steps. This results in improved productivity and cost reduction. Secondly, BPM provides organizations with better visibility and control over their processes, allowing them to monitor performance in real-time and make data-driven decisions[3]. This promotes timely interventions and continuous improvement. Thirdly, BPM fosters collaboration and coordination among various departments and stakeholders, facilitating effective communication and alignment of objectives. This leads to enhanced teamwork, quicker decision-making, and heightened customer satisfaction. Moreover, BPM empowers organizations to adapt and respond swiftly to evolving market conditions and customer needs, ensuring flexibility and a competitive edge. Overall, the development of BPM empowers organizations to achieve excellence in their processes, drive operational effectiveness, and foster sustainable growth [3], [5].

### B. Real-Time Enterprise

A Real-Time Enterprise refers to an organizational paradigm where operations are optimized to enable instant responsiveness and agility. This strategic approach involves harnessing advanced technologies and innovative methodologies to enhance business processes and decision-making capabilities[6]. By embracing the concept of a Real-Time Enterprise, companies can leverage real-time data access, proactive decision-making, and seamless collaboration to gain a competitive advantage in the market. This transformation requires integrating cutting-edge data analytics, real-time monitoring systems, and automated workflows to enable swift responses to market changes, customer demands, and emerging opportunities. Shifting towards a Real-Time Enterprise involves transitioning from traditional, time-consuming processes to agile methodologies, dynamic process modeling, and adaptive strategies driven by real-time insights. By embracing the Real-Time Enterprise vision, organizations position themselves for sustained growth, improved operational performance, and the ability to swiftly adapt to evolving market dynamics [7], [8].

### C. Time and Real-time Figures

The concept of real-time encompasses a wide range of meanings and finds extensive applications across various fields.

Its interpretation and utilization vary significantly, reflecting the diverse contexts and requirements in which it is applied. Real-time can refer to the ability to process and respond to data or events instantly, enabling rapid decision-making and actions. This is particularly crucial in time-sensitive industries such as finance, healthcare, and transportation. Additionally, real-time can also denote the synchronization of processes and activities with the passage of time, ensuring smooth coordination and minimizing delays [9], [10].

Real-time systems play a vital role in industries like manufacturing and logistics, where their reliance on such systems is substantial for enhancing operational efficiency and productivity [8]. These industries leverage real-time capabilities to optimize their processes, ensuring smooth and timely execution of tasks[5], [2]. For instance, the processing of large volumes of data in real-time enables the detection of anomalies and deviations, providing valuable insights for proactive decision-making and risk mitigation. This allows businesses to identify and address potential issues promptly, leading to improved operational performance and overall effectiveness. Therefore, the integration of real-time systems, coupled with advanced data processing techniques, proves instrumental in driving operational excellence across various sectors, including manufacturing and logistics [11].

Hence, the concept of real-time spans a vast spectrum of meanings and holds significant relevance in numerous sectors, highlighting its broad applicability and importance in today's dynamic and interconnected world.

Within the context of enterprise architecture development, the As-Is phase depicts the current state of the organization, while the To-Be phase represents the envisioned future state. This differentiation allows for a clear understanding of the present and future status of the company. In Business Process Model and Notation (BPMN), time constraints are effectively handled and modeled through the utilization of event time entities, which prove to be generally adequate for managing temporal aspects. However, as the demands of companies continue to expand across various dimensions, the time axis, known for its sensitivity and significance, becomes increasingly critical[9], [10].

### III. PROPOSED APPROACH

#### A. Right-Real-time Ontology

The initial phase involved establishing a formal definition of the real-time concept. While in an ideal scenario, real-time enterprises would respond instantaneously to business requirements, it is acknowledged that achieving true real-time capabilities is challenging. Hence, the concept of "near real-time" is introduced. The primary focus is on minimizing the time latency between data storage and availability, aiming to provide decision-makers with relevant and timely information. Considering the inherent challenges of achieving real-time capabilities in their entirety, it is important to prioritize timely reactions. Therefore, the key lies in delivering the appropriate information to the designated individual at the opportune moment. This ensures that decision-makers receive the necessary insights precisely when they are most advantageous [12].

When dealing with time in a real-time process, it becomes necessary to consider both an acceptance interval and a theoretical-time. The acceptance interval refers to a predefined range that ensures customer satisfaction, while the theoretical-time represents the ideal duration that can be predicted using various prediction tools. These two elements play a crucial role in effectively managing time within a real-time process.[12]. Based on the findings from the aforementioned results, we can provide a formal definition of real-time as follows:

$$\text{Right-Real-time ontology = Time ontology +} \begin{cases} \text{Latency} \\ \text{Acceptance interval} \\ \text{Theoretical time} \end{cases}$$

The adapted version of the time ontology, as depicted in Fig. 3, incorporates the concept of real-time [13]. This ontology defines time based on three components: time element, linear/nonlinear, and absolute/relative. However, our aim was to introduce a novel time component that would provide us with a fresh understanding of time.

In this adapted version, a new component called real-time is added to the time ontology. Within this component, three additional sub-components are included: latency time, acceptance interval, and theoretical time. This definition presents a new perspective on time, moving beyond the conventional notions of periods and calendars, and instead focusing on its significance in addressing the needs of real-time enterprises.

Please note that the figure mentioned can be found in Fig. 3, and the adapted ontology incorporates the real-time concept proposed by Kirikova et al. [13].

The first attribute in our definition is latency, which is closely linked to the concept of real-time as we previously mentioned. Since attaining real-time in its entirety is not feasible, latency becomes a reliable indicator in defining real-time. It can be conceptualized as an interval. In order to align with client needs, we found it essential to introduce a new attribute that establishes a safe range of latency. Thus, we defined the second attribute as the acceptance time interval. Similar to latency, the acceptance interval is also defined as an interval, representing a safety range within which latency does not disrupt the process flow.

Theoretical time, on the other hand, serves as our projection into the future based on our results. The purpose of defining the theoretical time is to enable result comparison and determine the degree to which we deviate from the ideal outcome. Generally, theoretical time is regarded as the ideal result, while the upper bound of the acceptance interval represents the worst outcome. Consequently, the lower bound of the acceptance interval corresponds to the theoretical time.

#### B. Right-real-time Process

The conventional interpretation of a process refers to a sequence of actions undertaken to attain specific outcomes. However, the definition of a process varies across different domains, tailored to suit their respective requirements. For instance, within the industrial context, a process represents a series of steps involved in manufacturing products. Within numerous enterprises and organizations, the concept of a process surpasses the simplicity of its elementary definition.
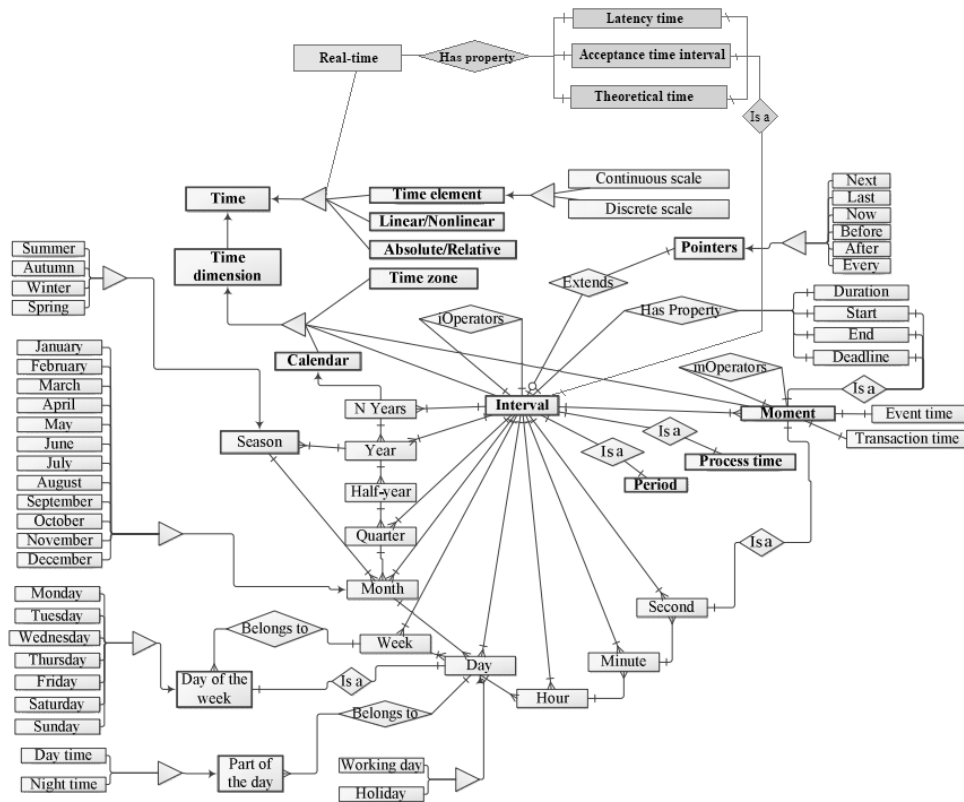
Fig. 3. Time ontology (adapted with real-time ontology).

It encompasses a greater level of complexity, often involving a collection of interconnected sub-processes, amplifying its intricacy. Companies have a wide array of languages at their disposal to effectively model processes, enabling them to visualize, interpret, and execute these processes while considering various constraints. One such language commonly utilized is BPMN (Business Process Model and Notation), which offers extensive capabilities in this regard.

Real-time processes encompass attributes that align with immediate responsiveness. However, the concept of real-time derives its definition from the specific requirements of customers, as opposed to a literal "right now" interpretation. So, if an event that happens an hour from now is judged acceptable, that occurrence is now practically the standard for what constitutes real-time. Establishing a precise definition for real-time is a complex task as it encompasses various interpretations. Given that achieving real-time in its absolute sense is often impractical, timely reactions become crucial. Therefore, delivering accurate information to the appropriate individual at the opportune moment becomes essential. Consequently, a real-time process refers to a process capable of providing a service and meeting client satisfaction. citeouarhim2019.

### C. Right-Real-time Enterprises/organisations

A real-time process constitutes a fundamental attribute of Real-Time Enterprises (RTE). These companies stand out because they can quickly react to different situations, usually using automated systems guided by built-in business rules or advanced technology solutions.

Within a real-time enterprise, there exist four primary categories of real-time processing, namely: (1) straight-through processing; (2) on-demand real-time data; (3) real-time performance management; and (4) real-time predictive analysis [14]. These processes encompass significant alterations over an extended duration, involving long-term changes, according to B. Kuglin and H. Thielmann [6]. These factors play a significant role in various areas: (1) within the internal and cross-company work processes; (2) in the division of labor both within an organization and across multiple companies; (3) through the implementation of technologies during the transition towards a Real-Time Enterprise; and (4) in the management of processes as well as the overall governance of the enterprise itself.

According to A. Ouarhim et al. [12] analysis, we can concisely outline the attributes of a real-time organization and establish a formal definition as follows: A real-time company is characterized by its agility, swift responses, prompt dissemination of information, rapid data analysis, efficient management of real-time processes, and incorporation of cutting-edge technologies, all of which converge to achieve near-zero latency.

### D. Managing Time in Business Process Toward Right Real-time Business Process in a Continuous improvement approach

*1) Time in a business process:* Within a business process, time manifests in various ways, yet a precise definition of the real-time concept remains elusive. In the context of a business process, time can be characterized as follows:

- Duration: interval with two ends (processes or tasks).

- Point: moments of execution of processes or execution time of tasks.

- Now: means that the process is auto-reactive (real-time).

- Time condition: reactive according to a time condition.

*2) Temporal latency index or temporal capability: real-time process validation algorithm:* The business process capability index serves as a metric that evaluates the correlation between a process's current performance and the predefined industry standards. It holds a keen interest in novel research about quality assurance and capability analysis. Capability indexes effectively measure both the potential and actual performance of processes, playing a vital role in quality improvement initiatives and serving as a cornerstone for successful implementation of quality programs [15] as shown in Fig. 4.

$$C_p = \frac{USL - LSL}{6\sigma}$$

$$C_a = 1 - \frac{|\mu - m|}{d}$$

$$C_{pk} = \min\left\{\frac{USL - \mu}{3\sigma}, \frac{\mu - LSL}{3\sigma}\right\}$$

Fig. 4. Capability indexes [15].

where:
• USL and LSL are the upper and the lower specification limits, respec-tively,
• Mu is the process mean,
• Sigma is the process standard deviation,
• m=(USL+LSL)/2 is the mid-point of the specification interval,
• d=(USL-LSL)/2 is half the length of the specification interval.

The process capability index, denoted as $C_p$, quantifies the overall variation of a process concerning the specified toler-ance, providing insight into the process's inherent potential or precision. On the other hand, the process capability index $C_a$ gauges the level of process centering, serving as an indicator when the process mean deviates from the target value, thereby reflecting process accuracy. Introducing the process capability index $C_{pk}$, it not only considers the magnitude of process variation but also accounts for the degree of process centering, thereby assessing process performance based on yield, which represents the proportion of conformity.

Undoubtedly, these measures serve as powerful tools for assessing process performance and efficiency. However, there exists another factor that significantly influences process effec-tiveness: time. Time has been a subject of ongoing research and, in this context, is increasingly recognized as a critical quality factor. By monitoring the behavior of time, valuable insights can be gained regarding the temporal dynamics of processes across different periods.

In light of this, we propose a comprehensive system for temporal process monitoring; see Algorithm 2. Each process is

characterized by its response time, waiting time, and execution time. Through this system, we aim to derive the following outcome [16]:

$$L_t = (T_{tmax} - T_{tmin})/(VAR * T_{tmax})$$

- Introducing the index $L_t$, also referred to as the "temporal latency index" or "temporal capability", provides valuable insights into the temporal dynam-ics of a process. This index enables us to gain a comprehensive understanding of a process's temporal behavior, empowering us to take timely action and implement necessary improvements accordingly.

- The parameter $T_t max$ represents the maximal tem-poral tolerance, which signifies the acceptable limit within each specific case study.

- $T_t min$ denotes the minimal temporal tolerance, rep-resenting the ideal scenario within the context of each case study.

- VaR: value at risk.

The concept revolves around determining the ratio between the tolerance margin and the maximum permissible risk of delay.

The percentage of response latency varies across different periods. This implies that if the error rate is calculated over a week, it tends to be higher compared to calculating it over a month. This variation can be attributed to the level of process discontinuity experienced within each period. Furthermore, the significance of response latency becomes less important when addressing past-present problems, whereas it holds greater importance when dealing with present-future problems. More-over, as organizations strive for real-time capabilities, acknowl-edging and mitigating time latency becomes imperative for achieving optimal results in dynamic environments. However, time latency plays a pivotal role in influencing the outcomes of various events and processes. The duration between the occurrence of an event and the corresponding system response can significantly impact the overall efficiency and effectiveness [17], [18].

Hence, it is crucial to determine the appropriate scale based on our specific requirements. If the problems or questions pertain to the present, our focus will be on the "day-month" scale. Conversely, if they relate to the future, our attention will be directed towards the "month-year" scale.

The utilization of VaR (Value-at-Risk) (see Algorithm 1) is driven by our interest in understanding the variations associated with time latency, which are never constant. By identifying the most unfavorable of these discrepancies, we can enhance control over our business processes. Hence, the index $L_t$ serves as a means to compare the theoretical and real outcomes. Value-at-Risk refers to the maximum potential loss that is only expected to occur with a given probability over a specific time period. In simpler terms, it represents the most severe loss anticipated within a defined time horizon, considering a certain level of trust [19].

Our approach entails identifying the most severe temporal risk that our business process can handle, utilizing a learning system. Therefore, it is crucial to determine the time period under examination, as previously discussed: "day," "week," or "month." When interpreting the VaR (value-at-risk) figure, one must consider the probability (x) and the holding period (t)[19]. First, we'll start by looking at each month of the year. After that, we'll analyze each month individually, and eventually, we'll broaden our examination to cover multiple years.

The utilization of the $L_t$ formula offers the advantage of simplifying the time period, making it applicable across all periods. Therefore, the choice of time period becomes significant in terms of result accuracy and precision.

---

**Algorithm 1** VaR calculation

---

**Result:** VaR value
i=1
**while** $i \leq lenght(filename)$ **do**
  tri=sort(Rnd(:,i))   **for** $k \leftarrow 1$ *to* $b(i) - 1$ **do**
  | $I(k) = k/(b(i) - 1)$
  **end**
  $J = find(I \geq 0.01)$   $PvaR = J(1)$   $VaRT(i) = -tri(PvaR)$a
**end**

---

**Algorithm 2** Real-time process validation algorithm

---

**Result:** Process state functioning
`// we start with algorithm inputs`
**Input:** $T_{tmax}$, $T_{tmin}$, VAR
$L_t = (T_{tmax} - T_{tmin})/(VAR * T_{tmax})$
**if** $0 < L_t \ll 1$ **then**
**else**
  | `// A requirement for making`
  |   `improvements in the process`
**end**
**if** $L_t \simeq 1$ **then**
**else**
  | `// A state of balance and indicates`
  |   `that the process is functioning`
  |   `well: low latency`
**end**
**if** $1 \ll L_t < 2$ **then**
**else**
  | `// The business process is operating`
  |   `at a near-perfect level:`
  |   `right-real-time process`
**end**

---

**Functioning of the index :**

The functioning of the index can be described based on the findings obtained from our analysis of three different case studies:

- If the value of $L_t$ deviates significantly from 1 to 0 (but never equals 0), it indicates the need for process improvements. Specifically, when $0 < L_t \ll 1$, action should be taken to enhance the business process.

- When the value of $L_t$ is approximately 1, it indicates a state of balance and indicates that the process is functioning well(on time demand: right-real-time): $L_t \simeq 1$.

- When the value of $L_t$ is significantly greater than 1 but less than 2 (and never equal to 2), it indicates that our business process is operating at a near-perfect level: $1 \ll L_t < 2$.

We have chosen to employ the historical method for calculating our VaR due to its simplicity, speed, and efficiency. By multiplying VaR with Ttmax, we obtain the maximum risk of daily delay. In other words, if tomorrow's delay is 'd', the worst-case scenario will be 'd + VaR*Ttmax'. Therefore, the concept of the $L_t$ formula represents a comparison between the maximum acceptable delay and the practical implementation using VaR (value at risk).

*3) Continuous temporal improvement approach:* Prevention is better than cure, having the appropriate tool to prevent specific issues will be better than reforming all present damage. Continuous improvement processes give many cycles that help to improve processes continuously, according to specific ethics for each cycle as we show in the following Table I, time in continuous improvement's tools and methods adapted version of [20]:

These continuous improvement processes have the goal of optimizing the performance of working conditions in terms of planning, organization, waste elimination, work methods, and knowledge management. In the present era, time wastage has emerged as a critical concern, particularly with the advent of the new digital transformation approach [5]. As significant changes continue to unfold, it becomes crucial to examine whether we can still maintain the same level of control over diverse business processes. Our pioneering contribution aims to tackle this challenge by not only preventing time wastage in business processes but also ensuring effective control to gain a comprehensive understanding of the real-time situation. Through our innovative approach, we strive to optimize time utilization and maintain a firm grip on business processes, thereby enabling informed decision-making and improved efficiency. First of all, our based contribution process is as follows: Fig. 5:

Our based process begins with analysis. We analyze data from our source of knowledge using simple and proactive parameters so we can detect the problem. After that, it is time to learn and discover different causes and try to find solutions to deliver the right information and knowledge to the right person at the right time. Before proposing solutions for the correction phase, we are faced with judging current practices.

We propose a rational approach dedicated to continuously improving time in business processes. As we all know, wasting time is a special case in all areas, and its damage becomes significant, especially when we face a serious situation. So, having a specific approach to time issues will give us the right answers to what we need. For that, we propose a right-real-time approach that is more responsive and compliant with time changes, namely our proposition about the $L_t$ index. Fig. 6, which is inspired by [21] and the Deming wheel introduced by William Edwards Deming (1950s).
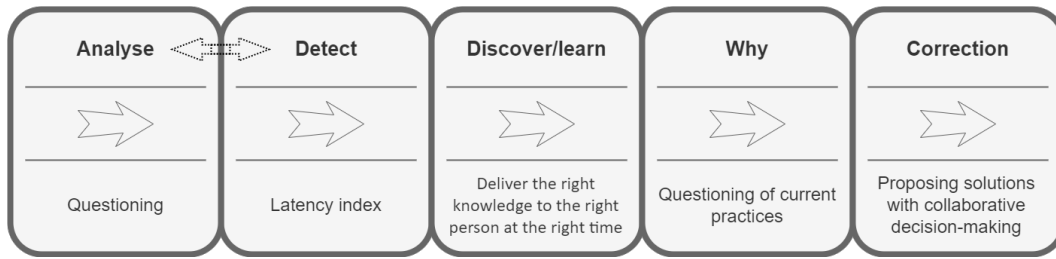
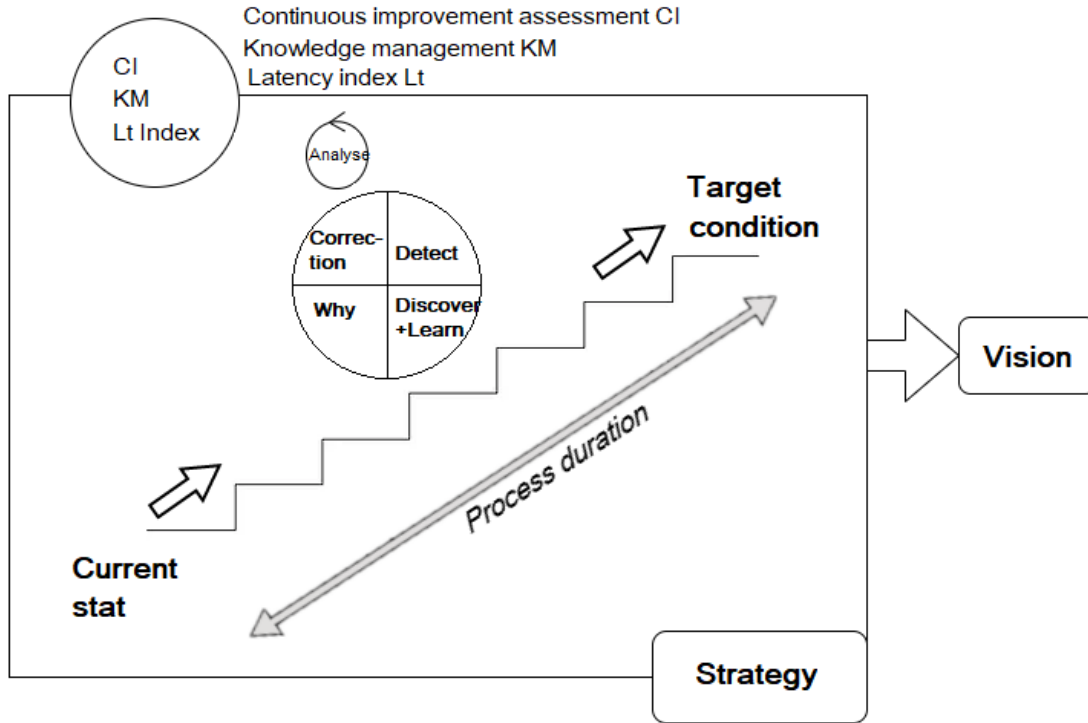Fig. 5. The based wheel process suitable to temporal issue, case right-real-time processes.



Fig. 6. Continuous temporal improvement approach and elements with Latency index Lt.

The originality of our approach can be seen in Fig. 6 as follows: The $L_t$ index provides real-time insight into the business process's temporal status. Based wheel process fits for right-real-time processes, serving continuous analysis or on-demand use for processes temporal improvements. In fact, during the implementation of continuous improvement, we must take into account the elimination of waste in all processes, highlighting the importance of our contribution to continuous improvement, which aims to control the waste of time to have the right-real-time CI process according to a previously established strategy. Our contribution $L_t$ index will support continuous improvement assessment and knowledge management as temporal knowledge that helps to evaluate the real-time situation of the process each time. Consequently, the company's vision will be clearer over time. The $L_t$ index will learn from all previous history of time processes to give a great prediction of present and future responding time process situations.

### E. Proposition of a New Component in BPMN: Right Real-time Component

Our focus lies on components that are time-related, and after an extensive analysis of the BPMN specification, we have identified Events as the key elements of interest. Events are directly related to FlowNodes and indirectly related to activities through BoundaryEvents and, specifically, CatchEvents. Both Events and Activities inherit from FlowNodes, which, in turn, inherit from FlowElements. The Process component inherits from FlowElementContainer, which is a composition of various FlowElements. This observation highlights that activities, along with the entire process, are more closely associated with time and real-time considerations, as indicated by the definition of time within a process [22].

Our proposal involves the creation of a novel component called Real-timeAttribute [16] as follows, Fig. 7, (OMG specification diagram [22]: adapted version): Fig. 7 illustrates a meta-model diagram depicting components that are relevant

TABLE I. OBSERVING TIME IN CONTINUOUS IMPROVEMENT'S
TOOLS/METHODS [ADAPTED TO [20]]

| Tool/method | Method Description | presence of waste-time analysis ($t$) |
|---|---|---|
| Kaizen Event | Kaizen events are structured initiatives that drive incremental improvements in processes, with a primary focus on enhancing process value and minimizing waste. | However, the connection is not straightforward. |
| Value Stream Map | A Value Stream Map is a visual depiction that illustrates the sequence and interrelationships of all the steps involved in a particular process. | No |
| Lead Time Analysis | The overall duration between the initiation and completion of a task, process, or service can be divided into two components: value-added time and non-value-added time. | Yes |
| Gemba | Japanese term used to describe the practice of physically going to the location where work is being performed. | No |
| 5 Why's | problem-solving method that involves repeatedly asking "why" to uncover the underlying root cause of a particular issue. | No |
| Spaghetti Diagram | A visual representation that illustrates the flow of transportation or movement of a product or service. | No |
| SIPOC | A comprehensive analysis of Supplier - Input - Process - Output - Customer that offers a customer-centric perspective of a process and its deliverable. | No |
| 6 S | An implementation method employed to establish and uphold a well-arranged work environment: Sort, Set in Order, Shine, Standardize, Sustain, and Safety. | However, the connection is not straightforward or direct. |
| Project Evaluation Matrix | A technique for assessing the business impact and the ease of resolving a problem to determine the priority of actions needed. | No |

to our study. The depicted figure highlights the focus on the rightRealTimeAttributes within the activity. Notably, this attribute has a direct impact on both Processes, Subprocesses, and Pools, as illustrated in the figure. As a result, these processes are automatically influenced by this particular property.

In Fig. 7, we can observe significant components that are relevant to our research. Additionally, we introduce our new component called Real-timeAttributes, which encompasses three essential attributes: latency, acceptanceInterval, and theoreticalTime (represented by RT). The inclusion of this component enhances our understanding and control over real-time processes, providing a clearer perspective and definition of what constitutes a real-time process.

The Fig. 7, represents a prototype of the extending component. We propose the development of a prototype that extends the Activity component, creating a specialized component known as "real-time Activity." This prototype aims to enhance the capabilities of the Activity component by incorporating real-time functionality and features.

## IV.   ARCHITECTURAL THINKING

### A.  Research Approach Diagram

Fig. 8, shows a diagram that elucidates the interconnection among all sub-sections and their corresponding outcomes, using the proposed based wheel process.

### B.  Capability Metamodel with our New Real-time Contribution

In the TOGAF content model, the objective of the organization is essential to fulfilling the capability. The TOGAF model can be further extended by providing additional meta-entities that describe the definition of capabilities as a measurable entity as shown in [23]. A business process enables the capability to execute the expected activities and outcomes. These entities that enable the capabilities, namely process, business service, and the lower level system components namely application architecture components, are measurable. A measurable entity is an entity whose attributes are measurable.

Our approach is to provide a new measure of capability concerning time. As the following figure shows inspired from [23], Fig. 9, our entity "Right-Real-time" is a measurable entity, that influences capability somehow according to each case study, for example, the influence can be generally latency as we discuss in the previous section. We define the measure of "Right-Real-time" as another entity "Right-Real-time-Index" that provides many attributes as shown in the previous section. "Right-Real-time-Index" has a goal that indicates the temporal situation of the process, as well we can make a decision, that we have named the index values'; that's considered as one of its attributes; previously as temporal latency index or temporal capability.

### C.  The Overview of the Unified Business Process Meta-model with our New Component: Real-time Business Process

So, we propose a unified business process meta-model containing our new component by integrating it into the unified meta-model proposed by Heidari, Farideh, et al. [24]. Their approach was to create a unified meta-model as a unified business process meta-model that provides a language-independent business process ontology. The mainstream business process modeling languages on which they were based are Business Process Modeling Notation (BPMN), Role Activity Diagram (RAD), Unified Modeling Language Activity Diagram (UML-AD), Integrated Definition for Function Modeling (IDEF0 and IDEF3), Structured Analysis and Design Technique (SADT), and Event-driven Process Chain (EPC). Each concept of these business process modeling languages is mapped onto only one concept in the unified business process meta-meta-model. They categorized the concepts of the unified business process meta-model into four aspects of a business process, namely: behavioral, functional, organizational, and informational aspects. This approach will give a full definition of the business process meta-model in terms of a unified meta-model and a real-time definition. Our CI approach plus real-time attributes, will give a specific time recognition about each process, which helps to first of all have a clear idea about the as-is timing situation and to make the right decisions about present and future processes. Fig. 10 presents our extended version of the unified modeling language; our components take part in behavioral aspects.
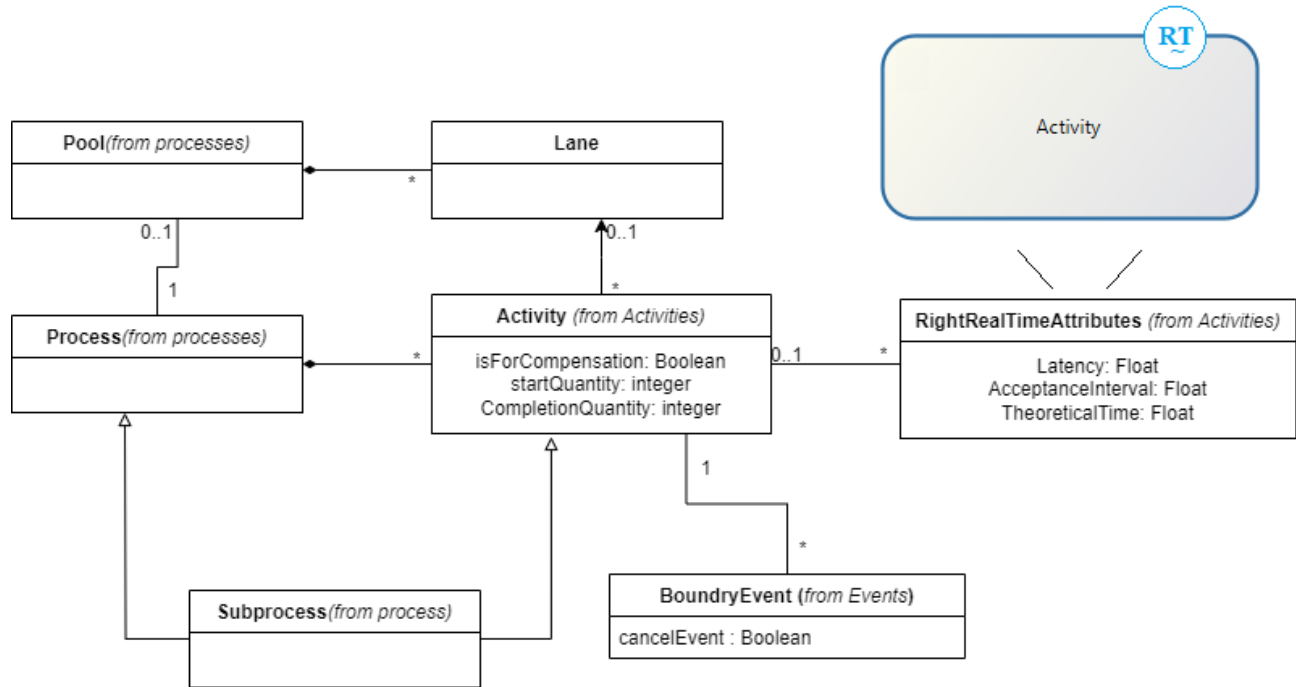
Fig. 7. New real-time component with additive attributes and a prototype of the extending component (RT).

## V. CASE STUDY: ANALYSIS, RESULTS AND DISCUSSION

### A. General Train Process

The train transport process is a crucial component of the modern transportation industry. It encompasses a series of coordinated activities that facilitate the movement of goods and passengers efficiently and safely. Starting from the scheduling and planning of train routes to ticketing, boarding, and on-board services, this process requires precision and attention to detail. Safety measures, maintenance routines, and adherence to schedules are fundamental aspects of this process. The efficiency of this process is vital for ensuring the seamless operation of train services, contributing to the overall connectivity and accessibility of regions and nations. Fig. 11 shows our proposed general train transport business process inspired by [25].

### B. Context of Work: Train Delays Problem

Train delay problems are due to so many reasons, especially train driver scheduling problems, which are considered more complex than other public transport problems. Indeed, this is due to driver work rules, constraints on the network, and the rolling stock. However, late trains can be resumed by: engineering work areas where the speed of trains is limited; the lack of double lines; conditioned speed (not the same all the way); overcrowded tracks owing to more and more trains each year; poor infrastructure bringing frequent maintenance, especially old tracks, which causes speed restrictions and delays; train driver behavior; and freight traffic contesting passenger routes.

According to Toor and Ogunlana [26], a delay is a result of many problems that can be resumed in factors related to local and environment, factors related to employees (designers, contractors, and consultants) and clients, and factors related to logistics sides such as lack of resources and other tasks problems such as planning and scheduling deficiencies. This problem is common in developed and developing countries and is considered one of the most recurring problems in construction projects. Problems of delay concern all types of construction projects, including trains. Major problems which this construction faces are usually due to three factors: system, resources, and communication.

However, the lateness of trains is due to so many reasons, like engineering work areas where the speed of trains is limited, lack of double lines, conditioned speed (not same all the way), Overcrowded tracks owing to more and more trains each year; poor infrastructure brings frequent maintenance, especially old tracks, which causes speed restrictions and delays; train driver behavior; and freight traffic contesting passenger routes. Trains driver scheduling problems are considered more complex than other public transport, and this is due to driver work rules, constraints on the network, and rolling stock.

This complexity arises from strict rules that the driver must follow for the safety of passengers and freight trains. Ronald et al. [27] did a study about the behavior and the psychological thinking of the train driver and set all situations that can make him not aware of his environment as a problem of control, and they present many methods that trait this kind of problems as COCOMO that help to understand train driver behavior.
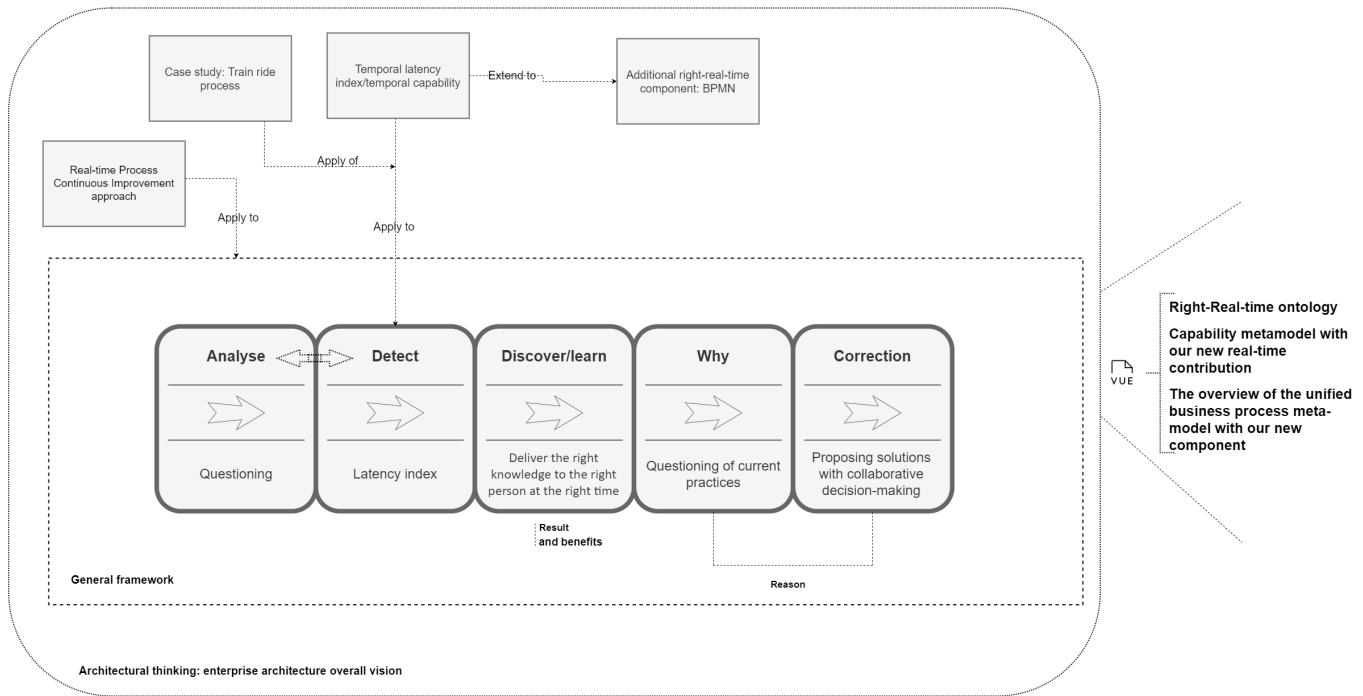
Fig. 8. Research approach diagram.

## C. Results

According to our based contribution process (see Fig. 5), we start by analyzing our data, then go to check the latency index value, and so on.

*1) Analyse:* As we have seen, the lateness of trains is due to so many reasons, like engineering work areas where the speed of trains is limited, the lack of double lines, conditioned speed (not the same all the way), overcrowded tracks owing to more and more trains each year, poor infrastructure that brings frequent maintenance, especially old tracks, which cause speed restrictions and delays, train driver behavior, and freight traffic contesting passenger routes. After analyzing our data based on the study, which focuses on presenting the findings specifically from year X, the results show that delays were not the same across all periods. Fig. 12, 13, and 14 show some examples of diagrams that we had during our analysis:

Fig. 12 shows delays of all trains/lines on each month during a year. Trains' late varies monthly and daily. As we can see, delays are not the same in all periods, but we can conclude that trains have the same attitude in all months. Fig. 13 shows the variation of mean delay in each station for the same train "train A" in a year, we see that delays are not the same and change according to each station.

Similarly, in the case of train B, as depicted in Fig. 14 during the same time-frame, it is observed that the mean delay varies across different stations. It's essential to note that in this context, the term "trains" serves as a representation of lines.

Our approach aims to identify any potential latency issues within our business process. In the event of a positive indication, we will proceed to explore and implement possible solutions. The advantage of using $L_t$ Index is to have a gain

TABLE II. TABLE OF RESULTS

| Years | $L_t$ |
|---|---|
| X | 0.142[1] |

[1] the value of $L_t$ deviates significantly from 1 to 0 (but never equals 0), it indicates the need for process improvements.

in terms of time. Indeed, it resolves the problem of latency twice.

*2) Detect:* Through extensive research conducted over various periods, this study focuses on presenting the findings specifically from the year 2018, offering a comprehensive overview. When implementing this methodology for the first time, it is advantageous to initially analyze previous years as a foundational element [13], [19], In line with this approach, we have specifically chosen to focus on previous years in our study. This approach provides a holistic perspective on the alignment of our processes with the concept of real-time. To facilitate a thorough examination, we gradually narrow down the time periods, starting with months and subsequently delving into weeks, and so forth [19]. This progressive analysis enables a deeper understanding and evaluation of our processes at different levels of granularity.

The company has already determined the values of $T_{tmax}$ and $T_{tmin}$ as two predefined elements. The analysis results will be condensed and presented in the following Table II.

Fig. 15 shows Train transport cases business process with our right real-time component.

*3) Discover/learn:* $L_t$ value shows that the train process is far from being a real-time business process.
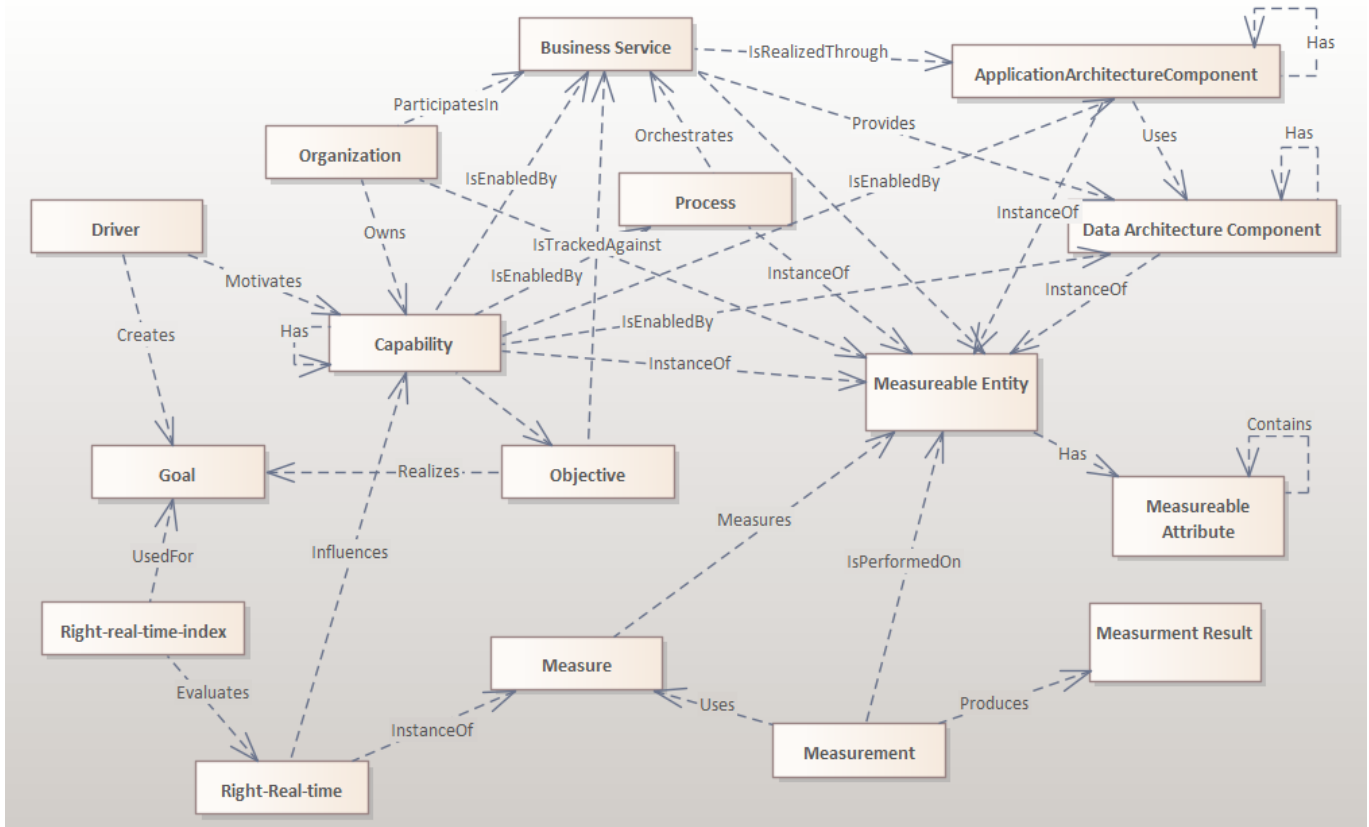We observed similar graphs for other trains and found that the

Fig. 9. Capability metamodel with right-real-time components (adapted with "right-real-time" components).
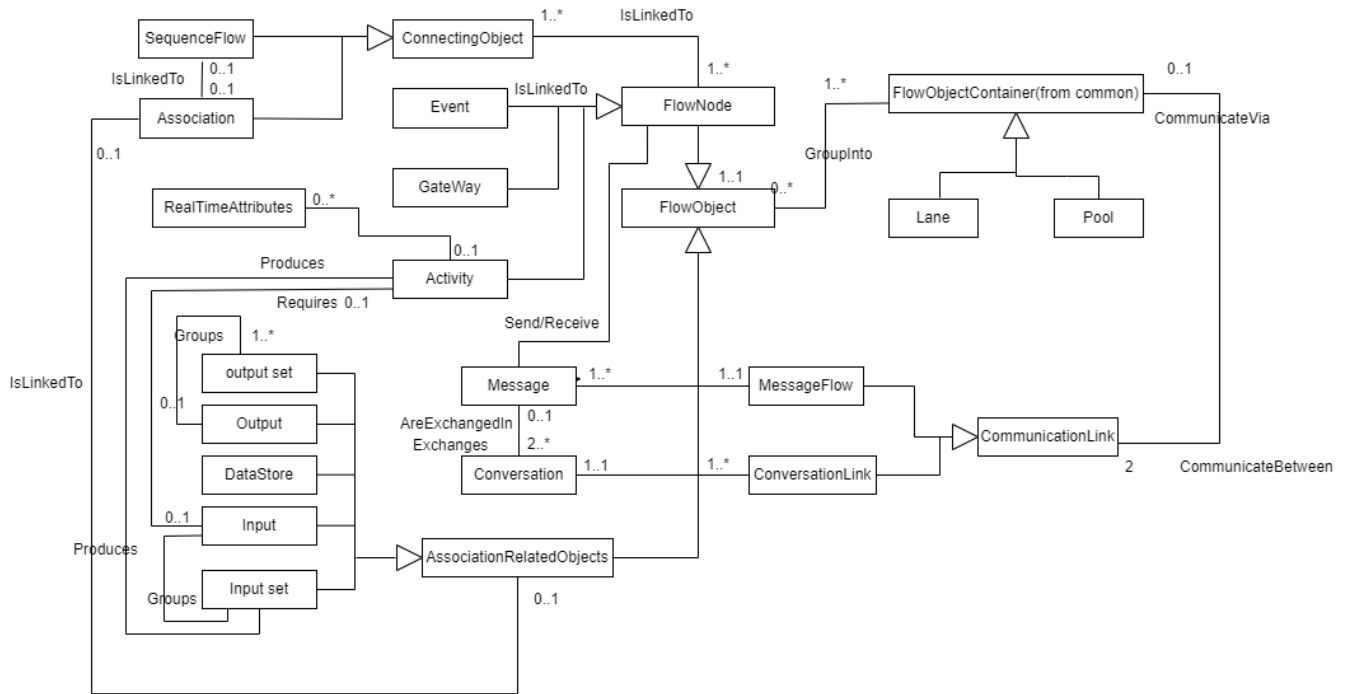


Fig. 10. The overview of the business process meta-meta-model with our new component.

variation in delays is related to periods in a year and station characteristics. The journey of each train is characterized by its lines and stations. A train can change lines daily. So, major delays along a line are related to busy stations compared to
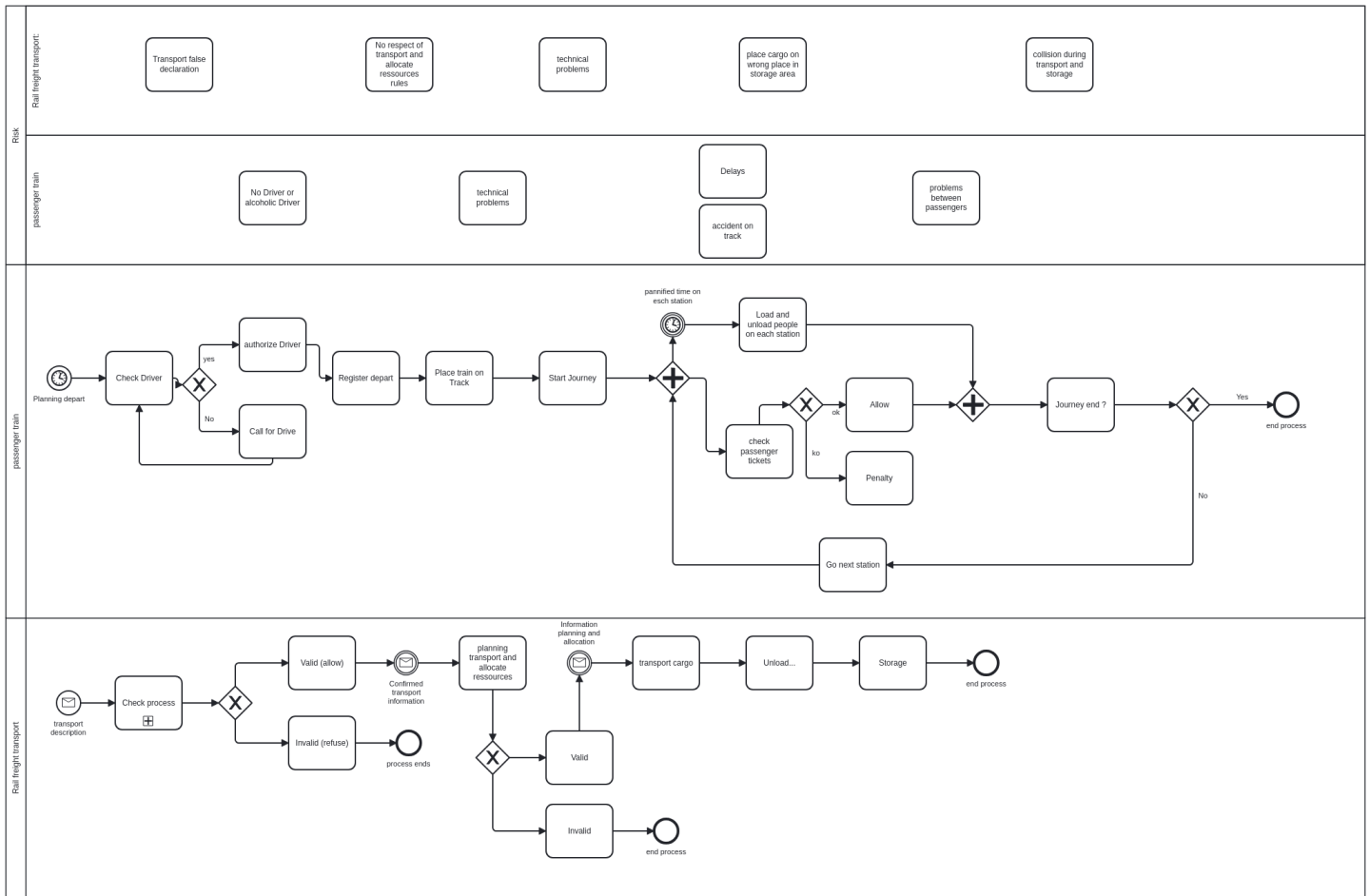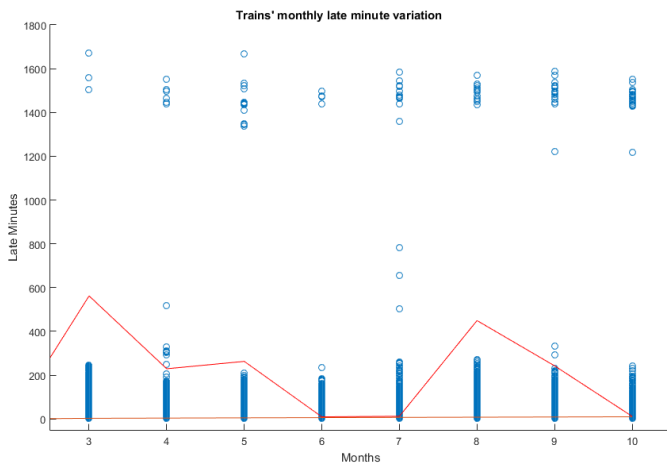
Fig. 11. Train transport business process.



Fig. 12. Trains' monthly late variation in the year X.



Fig. 13. Mean late minutes during train A's journey in the year X.

other lines with less busy stations.

*4) Why:* This problem of latency can be related to the variation in the load of people at each station and the need or demand for train transport during the year. There are periods when trains are less in demand in some lines; however, in

the same period, trains are more needed in other lines, which unfortunately causes delays.

*5) Correct:* In this phase, we identify and rectify errors, inconsistencies, or inefficiencies. This phase aims to refine and enhance the process's quality and performance by addressing any issues that have been identified during the assessment or execution stages. By conducting a thorough analysis, addressing issues, and optimizing, the corrective phase ensures
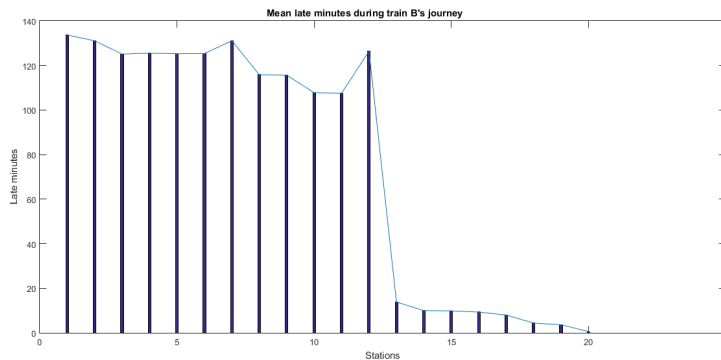
Fig. 14. Mean late minutes during train B's journey in the year X.

the alignment of the process with desired outcomes and adherence to established standards. This, in turn, enhances overall effectiveness and results. After the 'Why' phase, it is clear that a lack of stations and train traffic management is the main cause of latency. So, proposing a preliminary approach to managing train traffic according to the needs of each station will be the first step in improving this case study. The results attained will serve as inputs for the subsequent phase of our ongoing enhancement journey, creating a cycle of progress that perpetuates sequentially (see Fig. 6).

### D. Discussion

This approach combines two advantages: the first is the proposed based wheel process suitable for temporal issues, and the second is the latency index that indicates latency issues in a process. The usual waste tools didn't give a whole treatment that combined two approaches toward continuous temporal improvement.

## VI. Conclusion

We present a new "real-time process continuous improvement methodology" plus a "real-time process validation algorithm", which is an original search in terms of definition, modeling, and application. This approach is useful for every business process, including the time aspect, which allows them to identify and rectify any potential deviations, bottlenecks, or errors as they occur, preventing negative impacts on overall operations in terms of time. Indeed, time waste has emerged as a significant challenge in today's highly competitive world. Consequently, effectively managing and controlling time has become a critical endeavor. To address this issue, we have introduced a novel definition of the real-time concept that aligns with customer needs and our objectives of achieving a successful business process. Expanding on this definition, we've created an algorithm for real-time validation of business processes. The main goal is to ensure high-quality timing, with the ultimate aim of achieving minimal time latency (i.e., time latency $\simeq 0$). In essence, this algorithm serves as a means to assess process consistency by leveraging temporal capability. It provides decision-makers with insights into the temporal behavior of processes during execution, enabling them to make prompt decisions and find suitable solutions. The proposed algorithm is supported by our "continuous temporal improvement approach.". Furthermore, we have introduced a

new BPMN real-time component that includes various features to ease process monitoring within a continuous improvement (CI) approach. Furthermore, we have introduced a real-time unified business process meta-model that offers a comprehensive definition of the business process meta-model, unifying it with real-time considerations. By adopting our approach, organizations can gain specific insights into the temporal aspects of each process, establishing a clear understanding of the current timing situation and facilitating informed decision-making for both present and future processes. The limitation of the proposed approach lies regarding prediction; until now, we could use past data to evaluate the current state of a business process or how it could be if we didn't interact. So, in terms of perspective and future research, incorporating a deep learning tool into our approach would be advantageous for obtaining results.

## References

[1] Roeglinger M. Lehnert M., Linhart A. Exploring the intersection of business process improvement and bpm capability development: A research agenda. *Business Process Management Journal.*, 2017.

[2] Queiroz M. M.; Fosso Wamba S.; Machado M. C.; Telles R. Smart production systems drivers for business process management improvement: An integrative framework. *Business Process Management Journal.*, 26.5:1075–1092., 2020.

[3] P. Bazan and E. Estevez. Industry 4.0 and business process management: state of the art and new challenges. *Business Process Management Journal*, 28:62–80, 2022.

[4] Bhuiyan Nadia and Amit Baghel. An overview of continuous improvement: from the past to the present. *Management decision*, 43.5:761–771., 2005.

[5] Biernikowicz A. Gabryelczyk R., Sipior J. C. Motivations to adopt bpm in view of digital transformation. *Information Systems Management*, pages 1–17., 2023.

[6] Kuglin B. and H. Thielmann. The practical real-time enterprise: Facts and perspectives. *Springer Science and Business Media.*, 2005.

[7] et al. Ananyin, Vladimir I. Real-time enterprise management in the digitalization era. 13.1:7–17., 2019.

[8] Steiner W. Kopetz H. Real-time systems: design principles for distributed embedded applications. *Springer Nature*, 2022.

[9] Fuyuki Ishikawa Watahiki Kenji and Kunihiko Hiraishi. Formal verification of business processes with temporal and resource constraints. *IEEE international conference on systems, man, and cybernetics.*, pages 1173–1180., 2011.

[10] et al. Hashmi, Mustafa. Are we done with business process compliance: state of the art and challenges ahead. *Knowledge and Information Systems*, 57.1:79–133., 2018.

[11] Habeeb R. A. A.; Nasaruddin F.; Gani A.; Hashem I. A. T.; Ahmed E.; Imran M. Real-time big data processing for anomaly detection: A survey. *International Journal of Information Management*, 45:289–307., 2019.

[12] A. Ouarhim and K. Baïna. Towards a real-time business processes validation algorithm. *Procedia computer science*, 148:580–589., 2019.

[13] Ludmila Penicina Kirikova Marite and Andrejs Gaidukovs. Ontology based linkage between enterprise architecture, processes, and time. *New Trends in Databases and Information Systems: ADBIS 2015 Short Papers and Workshops, BigDap, DCSA, GID, MEBIS, OAIS, SW4CH, WISARD, Poitiers, France. Proceedings. Springer International Publishing*, 8-11:382–391, September. 2015.

[14] Gianmario Motta Barroero Thiago and Giovanni Pignatelli. Business capabilities centric enterprise architecture. *Enterprise Architecture, Integration and Interoperability: IFIP TC 5 International Conference, EAI2N 2010, Held as Part of WCC 2010, Brisbane, Australia. Proceedings. Springer Berlin Heidelberg.*, 20-23:32–43, September. 2010.

[15] W. L. Pearn Chen K. S. and P. C. Lin. Capability measures for processes with multiple characteristics. *Quality and Reliability Engineering International*, 19.2:101–110., 2003.
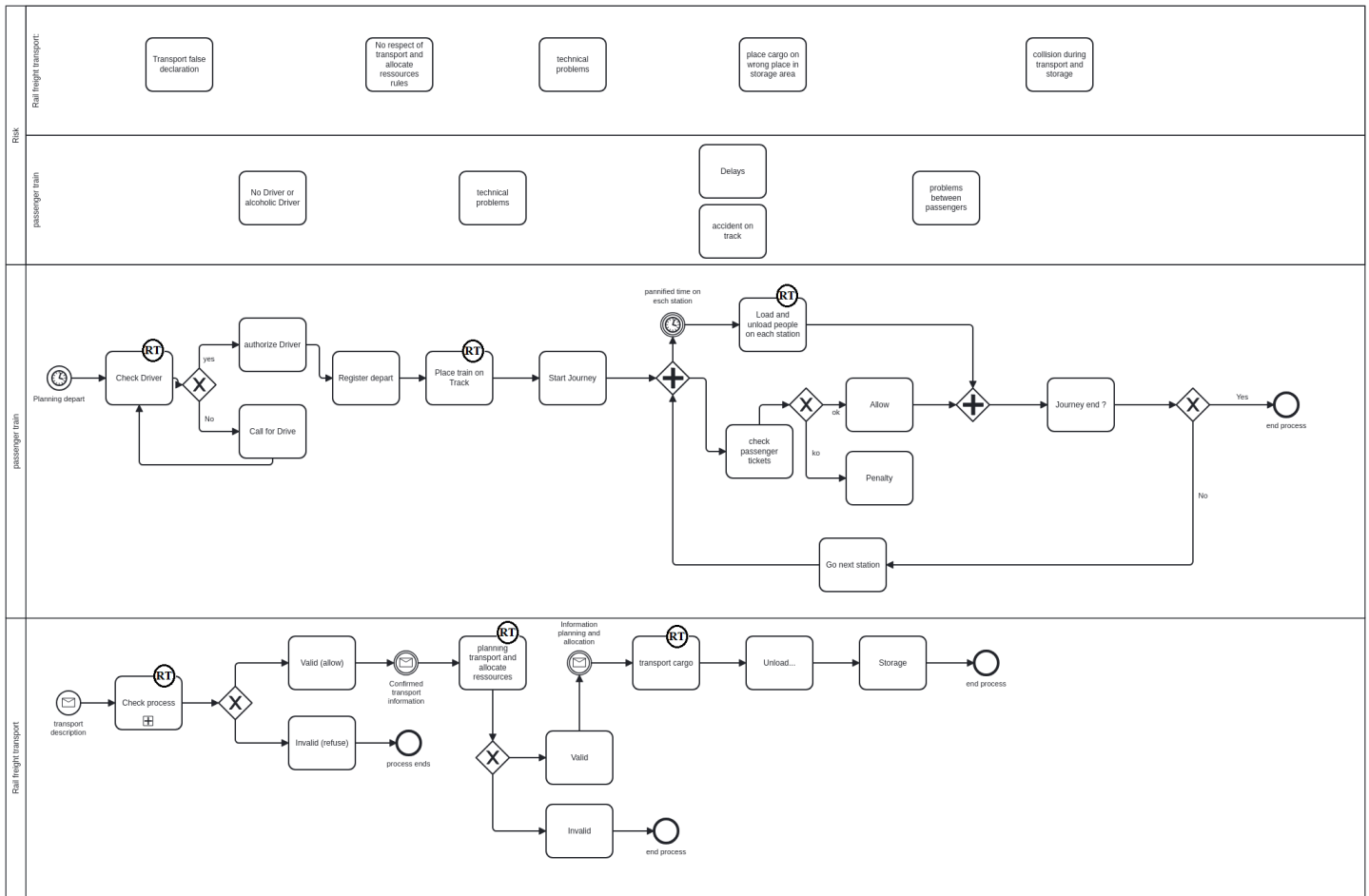
Fig. 15. Train transport business process with our right real-time component.

[16] Jihane Lakhrouit Ouarhim Asma and Karim Baïna. Business process modeling notation extension for real time handling-application to novel coronavirus (2019-ncov) management process. *5th International Conference on Cloud Computing and Artificial Intelligence: Technologies and Applications (CloudTech). IEEE*, pages 1–7., 2020.

[17] Ilona Boniwell and Philip G. Zimbardo. Balancing time perspective in pursuit of optimal functioning. *Positive psychology in practice: Promoting human flourishing in work, health, education, and everyday life*, pages 223–236., 2015.

[18] Terence R. Mitchell Morgeson, Frederick P. and Dong Liu. Event system theory: An event-oriented approach to the organizational sciences. *Academy of Management Review*, 40.4:515–537., 2015.

[19] Linsmeier Thomas J. and Neil D. Pearson. Value at risk. *Financial analysts journal*, 56.2:47–67., 2000.

[20] C W. C. Dicken. Bpm and lean: Part 1–the plan. 2012.

[21] M. Rother. Toyota kata: Managing people for improvement, adaptiveness and superior results. *MGH, New York*, 2019.

[22] OMG. Business process model and notation. *Object management group*, pages 1–7, 2013.

[23] Maureen. Du Toit, Francois A. et Tanner. A business architecture capability meta model and tool-set for providing function point estimation for enterprise architecture management. *Information Systems*, 650:4860, 2015.

[24] Heidari Farideh; Loucopoulos Pericles; Brazier Frances; et al. A meta-meta-model for seven business process modeling languages. *IEEE 15th Conference on Business Informatics.*, pages 216–221., 2013.

[25] et EL Fazziki Abdelaziz. Najib Mehdi, Boukachour Jaouad. A multi agent framework for risk management in container terminal: Suspect containers targeting. *Int. J. Comput. Sci. Appl.*, 10.2:33–52., 2013.

[26] Toor Shamas-Ur-Rehman and Stephen O. Ogunlana. Problems causing delays in major construction projects in thailand. *Construction management and economics*, 26.4:395–408., 2008.

[27] Guy H. Walker McLeod Ronald W. and Neville Moray. Analysing and modelling train driver performance. *Applied ergonomics*, 36.6:671–680., 2005.

# Modern Education: Advanced Prediction Techniques for Student Achievement Data

Xi LU

Hubei Institute of Fine Arts, Wuhan 430060, Hubei, China

*Abstract*—Enhancing educational outcomes across varied institutions like universities, schools, and training centers necessitates accurately predicting student performance. These systems aggregates the data from multiple sources—exam centers, virtual courses, registration departments, and e-learning platforms. Analyzing this complex and diverse educational data is a challenge, thus necessitating the application of machine learning techniques. Utilizing machine learning algorithms for dimensionality reduction simplifies intricate datasets, enabling more comprehensive analysis. Through machine learning, educational data is refined, uncovering valuable patterns and forecasts by simplifying complexities via feature selection and dimensionality reduction methods. This refinement significantly amplifies the efficacy of student performance prediction systems, empowering educators and institutions with data-driven insights and thereby enriching the overall educational landscape. In this particular research, the Decision Tree Classification (DTC) model is used for forecasting student performance. DTC stands out as a potent machine-learning method for classification purposes. Two optimization algorithms, namely the Fox Optimization (FO) and the Black Widow Optimization (BWO), are integrated to heighten the model's accuracy and efficiency further. The amalgamation of DTC with these pioneering optimization techniques underscores the study's dedication to harnessing the forefront of machine learning and bio-inspired algorithms, ensuring more precise and resilient predictions of student performance, ultimately culminating in improved educational outcomes. From the results garnered for G1 and G3, it is evident that the DTBW model demonstrated the most exceptional performance in both predicting and categorizing G1, achieving an Accuracy and Precision value of 93.7 percent. Conversely, the DTFO model emerged as the most precise predictor for G3, achieving an Accuracy and Precision of 93.4 and 93.5 percent, respectively, in the prediction task.

*Keywords—Student performance; classification; decision tree classification; fox optimization; black widow optimization*

## I. INTRODUCTION

The expansion of educational data sourced from admission systems, academic information systems, and e-learning platforms is substantial. Nonetheless, a significant portion of this data remains untapped due to its intricate nature and sheer volume. The analysis of this data holds pivotal importance in forecasting student performance. Data mining, known as knowledge discovery in databases (KDD), has proven to be efficacious across diverse domains, including education, paving the way for the emergence of Educational Data Mining (EDM) [1, 2].

Forecasting student outcomes in education significantly relies on EDM, allowing the anticipation of various results like passing, failing, and grading. A core focus involves establishing an early alert system to reduce costs, save time, and optimize available resources. Enhanced educational techniques are vital in refining student performance, enabling educators to tailor teaching methods and provide extra support where needed. These predictions empower students to gauge their potential academic progress and take necessary actions. Long-term institutional goals are centered on fortifying student retention, ultimately enhancing the institution's standing, rankings, and the career prospects of its graduates [3]–[6].

Educational establishments utilize data mining, commonly referred to as EDM, to thoroughly analyze the available data. Machine learning algorithms serve as pivotal tools for uncovering essential knowledge. Accurate performance prediction is instrumental in early identification of struggling students [7, 8]. EDM supports institutions in refining and developing novel learning methods by examining educational data. However, predicting academic performance presents challenges due to the diverse factors influencing it [9, 10]. Technological progress has facilitated the development of effective machine-learning methods [11–16]. Recent scholarly research emphasizes the efficacy of machine learning techniques in advancing the field of education.

Predicting student performance through machine learning (ML) is crucial for enhancing education in several ways. It enables early identification of academic struggles, allowing for timely interventions and personalized learning plans. By optimizing resource allocation and addressing factors influencing dropout rates, institutions can improve retention and graduation rates. Machine learning facilitates data-driven decision-making, adaptive assessments, and efficient educational planning. Continuous monitoring supports quality assurance, accountability, and a competitive advantage for institutions. Overall, it empowers educators to provide targeted support, leading to improved student outcomes and a more responsive education system.

## II. RELATED WORK

Ajay et al. [17] investigated the influence of the "CAT" social factor in predicting student performance among Indians. They employed four classifiers and found that the IB1 model exhibited the highest accuracy at 82%. This factor categorized individuals based on social status, directly impacting educational outcomes. Dorina et al. [18] developed a predictive model for student success using various classification algorithms. While the MLP model achieved the highest accuracy for identifying successful students, it encountered challenges in handling high-dimensional data and class

imbalances. Carlos used machine learning to create a student failure prediction model, achieving a high accuracy of 92.7% with the ICRM classifier. However, due to varying student characteristics, their study did not encompass testing across different educational levels. Edin Osmanbegovic et al. [19] devised a model to predict student academic success while tackling data dimensionality issues. Despite Naïve Bayes achieving the highest accuracy at 76.65%, the model did not effectively address the class imbalance problem.

A study [20] utilized various data mining methods to predict course dropouts in the context of EDM challenges. The support vector machine, with specific predictors, offered the most accurate classifications. However, including earned grades from prerequisite courses posed a limitation due to potential improvements in student knowledge during the course. Another study [21] aimed to enhance the ID3 model for predicting student academic performance, overcoming its inefficiencies in selecting attributes with numerous values. The proposed model significantly improved performance, achieving a high accuracy of 93% with the wID3 classifier. A study [22] introduced an early identification model for student failures, exploring multiple data mining methods and preprocessing techniques. Although the support vector machines outperformed other models, the study did not address reducing classification errors. Introducing an ensemble model, a study [23] aimed to identify underperforming students by combining classifiers. The ensemble model, incorporating standard-based grading assessments, outperformed individual classifiers, achieving an accuracy of 85%.

Suggesting a predictive system for online student learning performance, another study [24] found that methods considering time-dependent variables achieved higher accuracy. However, the model was not tested in an offline mode, potentially affecting its performance. Thammasiri et al. [25] proposed a model to predict poor academic performance among freshmen. The combination of support vector machines with SMOTE achieved the highest accuracy of 90.24%, addressing class imbalance issues. Challenging assumptions, a study [26] emphasized the applicability of data mining in small datasets for predicting student success. Although achieving over 90% accuracy with Reptree, the model did not effectively handle high data dimensionality or class balancing challenges. Addressing multiclass classification issues, a study [27] proposed a multi-level model to improve overall accuracy. This model, involving resampling and two levels of classification, achieved over 90% accuracy for both overall model and individual class predictions, using J48 as a key classifier.

## III. OBJECTIVE

The core aim of this research was to establish a robust machine-learning model designed for predicting Student Performance, drawing on data from credible sources. The study focused on leveraging the Decision Tree Classification (DTC) technique. An innovative approach was introduced by seamlessly integrating two optimization algorithms: Fox Optimization (FO) and Black Widow Optimization (BWO). This unique amalgamation of techniques was intended to significantly boost the Accuracy and Precision of the predictive model, thereby offering more effective forecasts of student performance within an educational setting. The DTC model is instrumental in predicting student performance in Mathematics due to its ability to comprehend and represent intricate relationships within data. Specifically tailored for educational contexts, the DTC method efficiently delineates critical factors influencing math performance. Its hierarchical structure allows for identifying significant decision paths, highlighting key determinants such as study habits, prior academic achievements, and socio-economic backgrounds. By comprehensively mapping these interdependencies, the DTC model predicts outcomes accurately and unveils pivotal insights essential for targeted interventions and tailored academic support, thereby enhancing student performance in Mathematics.

This study underscores the vital role of data-driven predictive models in education, advocating for a comprehensive approach to evaluate students' academic performance. Demonstrating the effectiveness of data mining techniques, including clustering and classification, the research innovatively integrates the DTC model with FO and BWO. This integration highlights the potential of combining machine learning and optimization algorithms to enhance precision, providing a robust toolkit for addressing challenges in students' academic journeys. The thorough evaluation process reveals the significant potential of these hybrid models to improve the DTC model's classification accuracy and precision, contributing to advancements in academic performance prediction.

## IV. MATERIALS AND METHODOLOGY

### A. Data Preparation

The primary aim of this study revolves around constructing a robust method to accurately evaluate students' academic performance while considering various contextual factors that influence it. To accomplish this objective, the initial dataset necessitates crucial preprocessing steps. The first essential step involves converting textual data into numerical values, a foundational requirement for conducting machine learning tasks. This conversion is pivotal as it facilitates effective data analysis and enables the application of advanced statistical techniques. The dataset encompasses a diverse range of variables that potentially impact students' academic outcomes, encompassing factors such as sex, school, urban or rural residency (address), age, family size (famsize), parental cohabitation status (Pstatus), parental education and occupations (Medu, Fedu, Mjob, and Fjob), school choice motivation (reason), weekly study time (studytime), guardian, home-to-school travel time (traveltime), current health status, past class failures (failures), participation in supplementary education (schoolsup), family educational support (famsup), engagement in extra paid classes, involvement in extracurricular activities, attendance at nursery school, aspirations for higher education, access to the internet, student absences, weekday (Dalc), and weekend (Walc) alcohol consumption, involvement in romantic relationships, quality of family relationships, free time, and frequency of socializing.

This research aims to predict and categorize students' academic performance, utilizing the G1 and G3 variables. G3 represents final grades obtained from school reports, ranging

from zero (indicating the lowest grade) to 20 (representing the highest grade). These grades are segmented into four distinct levels: Poor (0–12), Acceptable (12–14), Good (14–16), and Excellent (16–20), allowing for a more nuanced evaluation of student achievement. This methodology seeks to establish a comprehensive framework for comprehending and assessing academic performance within a myriad of contextual factors, ultimately contributing to improvements in educational practices and developing policies in the academic sphere.

Fig. 1 displays a correlation matrix detailing the relationships among input and output variables within this study. It notably highlights the positive influence of parental education, particularly maternal education, on students' academic performance. Moreover, factors such as daily and weekly alcohol consumption, prior academic failures, and student age demonstrate discernible impacts on school grades. Ultimately, the matrix underscores the critical importance of both study time and parental education as pivotal factors

contributing to academic success. Notably, there is a strong positive correlation (0.8264) between grades in the first period ("G1") and final grades ("G3"), indicating that students who perform well in the initial period tend to have higher final grades. Additionally, some demographic and lifestyle factors exhibit correlations. For instance, parental education levels ("Medu" and "Fedu") show moderate positive correlations, implying a potential influence on academic performance. The variable "sex" demonstrates a weak negative correlation with "age" (-0.0437), suggesting a slight tendency for younger students to be male.

The correlation matrix provides a snapshot of associations between different variables, offering insights into potential patterns and relationships. However, it is important to approach these correlations cautiously, as correlation does not imply causation and other factors may contribute to the observed relationships.
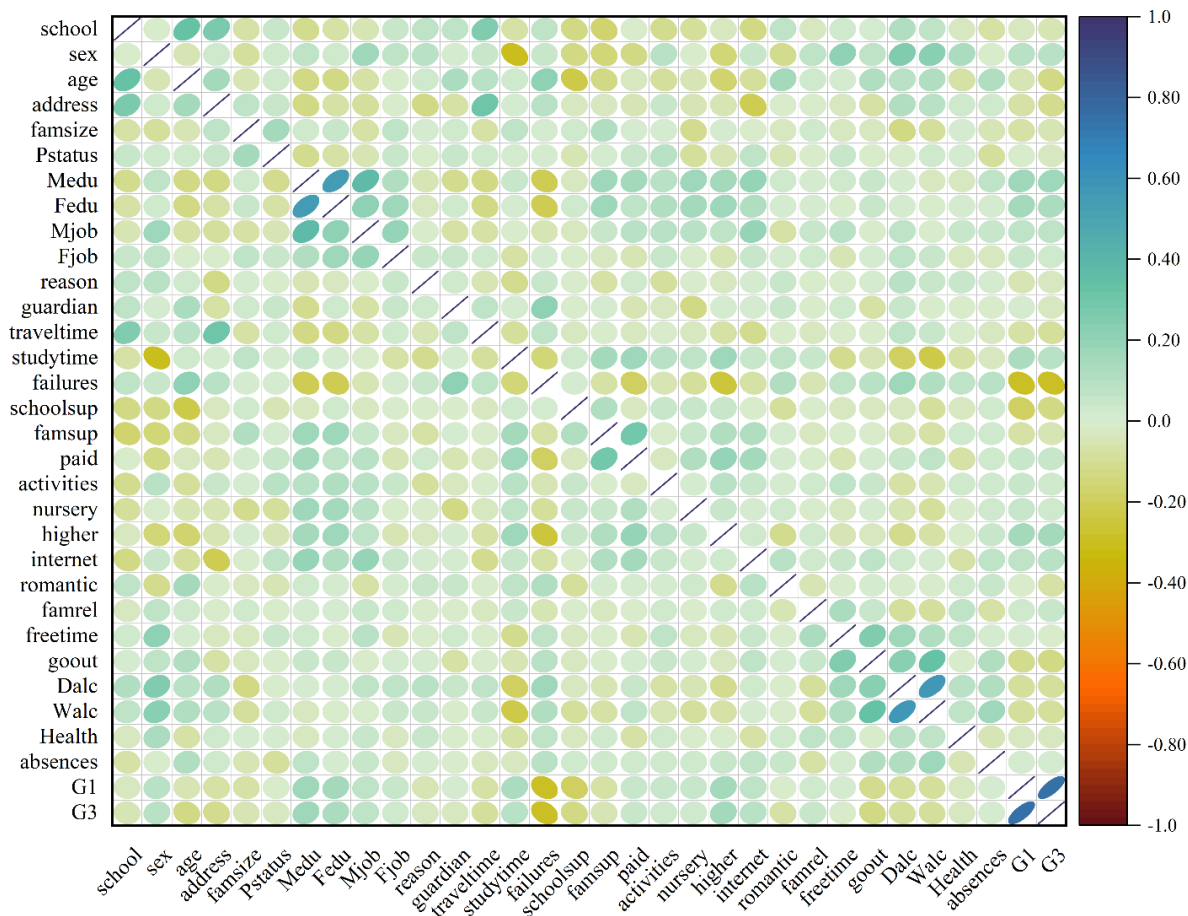


Fig. 1. Correlation matrix for the input and output variables.

## B. Evaluation of Models' Applicability

In academic studies focused on classification problems, Accuracy is a widely employed metric used to evaluate the overall performance of a model. It relies on four fundamental components: True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN). TP signifies accurate predictions, TN represents correct negative predictions, FP indicates incorrect positive predictions, and FN denotes

inaccurate negative predictions. However, Accuracy tends to favor the majority class, offering limited insights in situations where data is imbalanced. Three additional evaluation metrics—Recall, Precision, and F1-Score—are utilized to overcome this limitation. Recall evaluates the model's capability to correctly identify all relevant instances within a specific class, which is crucial in reducing False Negatives. Precision measures the accuracy of positive predictions, aiming

to minimize False Positives, instances predicted as positive but not belonging to the class. F1-Score, combining Precision and Recall, provides a balanced assessment of model performance, particularly valuable in scenarios with imbalanced data, considering both minority and majority classes. Defined by mathematical equations, these metrics collectively provide a deeper understanding of a classification model's effectiveness. They are especially beneficial in challenging situations involving imbalanced data, where the interpretation of Accuracy might be misleading. The utilization of these metrics empowers researchers and data analysts to make more informed decisions and adjustments to enhance model performance in such intricate scenarios [28].

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \qquad (1)$$

$$Precision = \frac{TP}{TP+FP} \qquad (2)$$

$$Recall = TPR = \frac{TP}{P} = \frac{TP}{TP+FN} \qquad (3)$$

$$F1\_score = \frac{2 \times Recall \times Precision}{Recall+Precision} \qquad (4)$$

### C. Decision Tree Classification (DTC)

A decision tree takes the form of a structure resembling a flowchart, where each internal node conducts a test based on an attribute, with each branch signaling the outcome of that particular test. Meanwhile, every leaf node, also termed a terminal node, denotes a distinct class label. Making predictions with a decision tree involves assessing the attribute values of a given data point, typically referred to as a tuple, by following a path from the root of the tree to a leaf node containing the projected class label for that specific data point. The strength of decision trees lies in their ease of conversion into classification rules. They serve as predictive models in decision tree learning, enabling the translation of observations about an object into conclusions about its intended value. These models have diverse applications in statistics, data mining, and machine learning, particularly in classification trees, which specifically handle finite class values. Compared to other classification methods, decision tree construction is commonly recognized as a swift process [29].

The decision tree relies on three key parameters:

*1) D (Data Partition):* D represents the initial dataset containing training examples and their respective class labels.

*2) Attribute list:* This parameter comprises attributes that detail the features of the data.

*3) Attribute selection method:* This parameter defines the strategy used to select the most suitable attribute for creating divisions or branches in the decision tree. Common methods involve measures like information gain or the Gini index.

Here is an overview of how the algorithm operates:

- It initiates by establishing a node labeled "A."

- If all the examples in the present dataset share the same class, "A" becomes a leaf node designated with that common class label.

- When the attribute list is empty, node "A" transforms into a leaf node, now tagged with the class that most frequently appears among the data samples.

- The algorithm then selects the attribute to split the data in a way that generates the purest subsets.

- Node "A" is assigned this selected attribute as the decision criterion.

- If the chosen attribute is discrete, it is removed from the attribute list.

- The data is segregated into subsets based on the outcomes of the selected attribute.

- If any of these subsets are empty, a leaf node is linked to node "A," labeled with the majority class of the original dataset.

- For non-empty subsets, the process repeats recursively, commencing with the creation of a new node until all data partitions have been addressed.

- Ultimately, the algorithm returns the resulting decision tree structure.

This algorithm is a foundational process for constructing decision trees, commonly applied in tasks involving data classification and predictive modeling within machine learning and data analysis contexts.

DTC is a preferred method for predicting student performance due to its interpretability, ability to handle non-linear relationships, versatility with mixed data types, ease of implementation, and avoidance of overfitting through pruning. DTC is suitable for educational datasets with both categorical and numerical variables, making it applicable to real-world scenarios. Additionally, decision trees can be part of ensemble methods, offering improved predictive accuracy. The transparency of decision tree models is valuable in educational contexts, enabling stakeholders to understand and discuss predictions.

### D. Fox optimization (FO)

The Fox Optimization Algorithm (FO) draws inspiration from the hunting behavior of red foxes and is structured around two primary phases: exploitation and exploration. The exploitation phase mimics a fox closing in on its prey, utilizing strategies to optimize the immediate vicinity. Conversely, the exploration phase is influenced by the relative distance between the fox and its target. This algorithm functions with a consistent population of foxes, maintaining a set structure as detailed below [30]:

$$\bar{x} = (x_0, x_1, \dots, x_{n-1}) \qquad (5)$$

In the identification of each fox $\bar{x}^t$ within the t-th iteration, a notation $\left(\bar{x}_j^i\right)^t$ is introduced. In this context, $i$ represents the count of foxes, while j denotes the specific coordinates within the solution space, delineated by the dimensions. $(\bar{x})^{(i)} = \left[(x_0)^{(i)}, (x_1)^{(i)}, (x_2)^{(i)}, \dots, (x_{n-1})^{(i)}\right]$ is employed to denote each point within the solution space $< a, b >^n$, where $a, b \in \mathbb{R}$. Furthermore, with regard to the solution space, a function

$f \in \mathbb{R}^n$ is regarded as the standard function of n variables. If the value of this function, $f\left((\bar{x})^{(i)}\right)$, represents a global maximum or minimum within the interval $< a, b >$, then $\left((\bar{x})^{(i)}\right)$ is deemed the optimal solution.

When foxes struggle to find prey, family members embark on the quest for food. When a more promising area is discovered, they share and communicate this location within the population, effectively supporting it and considering the associated cost. The metric utilized for this dissemination relies on the Euclidean squared distance.

$$D((\bar{x}^i)^t, (\bar{x}^b)^t) = \sqrt{\|(\bar{x}^i)^t - (\bar{x}^b)^t\|}, \qquad (6)$$

$(\bar{x}^b)$ represents the individuals within the population shifting their positions towards the direction of the best performer.

$$(\bar{x}^i)^t = (\bar{x}^i)^t + \alpha * S * ((\bar{x}^b)^t - (\bar{x}^i)^t), \qquad (7)$$

Here, $\alpha$ is randomly chosen from the range $\left(0, d((\bar{x}^i)^t, (\bar{x}^b)^t)\right)$, while S signifies the 'sign' word. The random value $\beta$, ranging between 0 and 1, remains consistent for all individuals in the population. This value embodies the behavior of the fox as:

$$\begin{cases} Stay\ and\ masquerade & if\ \beta \le 0.75 \\ Move\ closer & if\ \beta > 0.75 \end{cases} \qquad (8)$$

An advanced Cochleoid equation elucidates the behavior of individuals when $\beta$ influences the movement of the population in a given iteration. Two components determine the fox radius: $\phi_0 \in < 0,2\pi >$ representing the initial observation angle, and $\alpha \in < 0,0.2 >$ as a scaling parameter. This value is preset for all individuals in the population, symbolizing random alterations in distance as the fox approaches the target.

$$r = \begin{cases} a\,\frac{sin\phi_0}{\phi_0} & if\ \phi_0 \ne 0 \\ \delta & if\ \phi_0 = 0 \end{cases} \qquad (9)$$

$$\begin{cases} x_0^{new} = ar * cos(\phi_1) + x_0^{ac} \\ x_1^{new} = ar * sin(\phi_1) + ar * cos(\phi_2) + x_1^{ac} \\ x_2^{new} = ar * sin(\phi_1) + ar * sin(\phi_2) + ar * cos(\phi_3) + x_2^{ac} \\ \quad\quad\quad ... \\ x_{n-2}^{new} = ar * \sum_{q=1}^{n-2} sin(\phi_q) + ar * cos(\phi_{n-1}) + x_{n-2}^{ac} \\ x_{n-1}^{new} = ar * sin(\phi_1) + ar * cos(\phi_2) + \cdots + ar * sin(\phi_{n-1}) + x_{n-1}^{ac} \end{cases} \qquad (10)$$

In this context, $\delta$, fluctuating between 0 and 1, stands as a random value set at the beginning of the algorithm, contingent upon prevailing weather conditions. The movement pattern for the population of individuals is articulated as follows:

Where "ac" in $x_0^{ac}$ signifies "actual," and $\phi_1, \phi_2, \phi_3$, and so on, up to $\phi_{n-1}$, all exist within the range of $< 0,2\pi >$.

5% of the least successful candidates are selected based on the criterion function to replicate this action in each iteration. This selection is a subjective assumption aimed at introducing slight variations within the group. In iteration t, the two top-performing individuals are chosen for an alpha couple.

The pair comprises $(\bar{x}^{(1)})^t$ & $(\bar{x}^{(2)})^t$, while the center of the habitat is calculated using a specific equation. The square of the individual Euclidean distance between the couple determines the habitat range.

$$(H^{cntr})^t = \frac{(\bar{x}^{(1)})^t + (\bar{x}^{(2)})^t}{2} \qquad (11)$$

$$(H^{diamtr})^t = \sqrt{\|(\bar{x}^{(1)})^t - (\bar{x}^{(2)})^t\|} \qquad (12)$$

In this context, 'H' denotes the Habitat. Each iteration involves the selection of a random parameter 'q' ranging from 0 to 1, governing the substitutions conducted throughout the repetition in the following manner:

$$\begin{cases} Reproduction\ Of\ The\ Alpha\ Couple \\ \quad\quad if\ q < 0.45 \\ New\ Nomadic\ Individual \\ \quad\quad if\ q \ge 0.45 \end{cases} \qquad (13)$$

The top two candidates indicated as $(\bar{x}^{(1)})^t$ and $(\bar{x}^{(2)})^t$, are amalgamated to generate a new candidate, denoted as $(\bar{x}^{(rep)})^t$, where "rep" signifies reproduction. This fusion takes place in the following manner:

$$(\bar{x}^{(rep)})^t = q\,\frac{(\bar{x}^{(1)})^t + (\bar{x}^{(2)})^t}{2} \qquad (14)$$

The Steps of the Fox Optimization algorithm is represented as Algorithm 1.

ALGORITHM. 1. PSEUDO-CODE OF FO

Commence,
Establish the algorithm's parameters: the fitness functions $f(0)$, the number of iterations $T$, the initial fox observation angle $\phi_0$, the maximum population size $n$, weather conditions $\theta$, and the solution space range $< a, b >$,
Create a population of $n$ foxes randomly distributed within the solution space.
t= 0
while $t \le T$ do
Define iteration coefficients: fox proximity change ($\alpha$), scaling parameter ($\alpha$).
For every fox within the current population,
Organize individuals based on their fitness function values,
Select $(\bar{x}^b)^t$
Compute the repositioning of individuals
If the new position is superior to the previous one, then
Relocate the fox to the new position,
else
Revert the fox to its previous location,
end if
Determine the parameter $\beta$ to define the fox's hunting awareness,

If the fox remains unnoticed, then
Calculate the fox's observation radius ($r$)
Compute the repositioning
else
The fox maintains its current position to remain concealed,
end if
end for
Arrange the population following the fitness function,
Eliminate the poorest-performing foxes from the group, or they fall victim to hunters,
Introduce new foxes into the population as nomadic foxes outside the habitat or through reproduction from the alpha couple within the herd
t + +,
end while
Return the fittest fox $(\bar{x})^b$,
Stop.

### E. Black Widow Optimization (BWO)

The BWO is a recent and intriguing meta-heuristic approach for tackling complex numerical optimization challenges [31]. BWO incorporates operators commonly found in evolutionary algorithms, akin to Genetic Algorithms (GAs) [31]. Like other evolutionary algorithms, BWO employs criteria resembling natural evolutionary processes, such as selection, reproduction, and mutation, which vary and distinguish it from other evolutionary methods. However, what sets BWO apart is its simulation of the unique mating behavior of black widow spiders. Furthermore, BWO exhibits distinctions from traditional evolutionary algorithms, contributing to its strong performance in solving complex problems. This algorithm draws inspiration from Darwin's theory of natural selection, characterized by species evolving and the emergence of new ones. BWO is known for its rapid convergence and ability to evade local optima, making it well-suited for solving various optimization problems with multiple local optima. This success is attributed to BWO's balanced approach, maintaining harmony between the exploration and exploitation phases. For a visual representation of the BWO process (see Fig. 2).
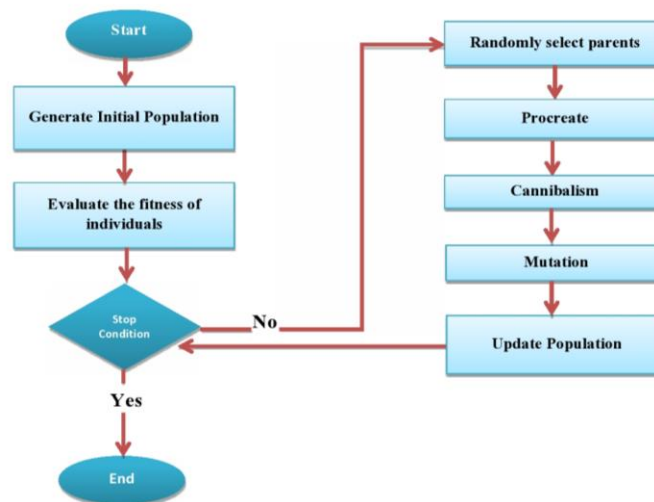


Fig. 2. Flowchart of the BWO.

The primary steps of the BWO can be summarized as follows:

#### 1) Step one: Initialization

During this step, the population consists of a specific number of widows, denoted as N, where each widow is represented as an array of size $1 \times N_{var}$, signifying a potential solution to the problem. This array can be described as follows: $widow = (x_1, x_2, \ldots, x_{N_{var}})$, where $N_{var}$ corresponds to the dimensionality of the optimization problem. $N_{var}$ can also be understood as the count of threshold values the algorithm aims to determine. Here, $x_i$ represents the $i - th$ candidate solution within the array.

The fitness of each widow is determined by evaluating the fitness function, denoted as f, for every widow in the set $(x_1, x_2, \ldots, x_{N_{var}})$. This fitness value can be expressed as follows: fitness = f(widow), which is equivalent to $fitness = f(x_1, x_2, \ldots, x_{N_{var}})$. The optimization procedure commences by initializing a population of spiders randomly in a matrix of dimensions $N_{pop} \times N_{var}$. Subsequently, pairs of parents are selected randomly to engage in the reproduction step, which is followed by the mating process. During or after mating, the male black widow is consumed by the female.

#### 2) Step two: Procreate

During the procreation step, an alpha (α) array is generated. This alpha array has the same length as a widow array and is filled with random numbers. Subsequently, offspring is

generated using alpha (α) and Eq. (14), where $x_1$ and $x_2$ represent the parents and $y_1$ and $y_2$ denote the offspring. The outcome of the crossover operation is assessed and then stored for further processing.

$$y_1 = \alpha \times x_1 + (1 - \alpha) \times x_2 \text{ and } y_2 \qquad (15)$$
$$= \alpha \times x_2 + (1 - \alpha) \times x_1$$

### 3) Step three: Cannibalism

The cannibalism process can be classified into various categories, including sexual cannibalism, sibling cannibalism, and a commonly observed form in which baby spiders consume their mother. Following the implementation of the cannibalism mechanism, the resulting new population is assessed and saved in a variable referred to as $pop2$.

### 4) Step four: Mutation

The mutation process involves randomly selecting a number of individuals, denoted as $Mutepop$, from the population to undergo mutation. Each selected solution has two elements within their array randomly exchanged in this mutation operation. After applying mutation, the resulting new population is evaluated and stored in a new population variable, typically named $pop3$. Finally, the new population is obtained by combining (or migrating) the individuals from $pop3$ and $pop2$. Subsequently, this combined population is sorted, aiming to identify the best widow with $N_{var}$ dimensions in terms of threshold values. Algorithm 2 provides the pseudo-code for the BWO algorithm.

---

ALGORITHM 2: PSEUDO-CODE OF BWO ALGORITHM

Initialize: Maximum number of iterations, rate of procreating, rate of Cannibalism, rate of mutation;
while **Stopconditionnotmet** do
for $i = 1$ to $nr$ do
Randomly select two solutions as parents from $pop1$.
Generate D children
Destroy father.

---

Based on the cannibalism rate, destroy some of the children (newly achieved solutions).
Save the remaining solutions into $pop2$.
end for
Based on the mutation rate, calculate the number of mutation children $nm$.
for $i = 1$ to $nr$ do
Select a solution from $pop1$.
Mutate randomly one chromosome of the solution and generate a new solution.
Save the new one into $pop2$.
end for
Update $pop = pop2 + pop3$.
Returning the best solution.
Return the best solution from pop.
end while

---

## V. RESULTS AND DISCUSSION

### A. Convergence Results

In this study, two powerful metaheuristic optimization algorithms, the FO and BWO, were employed to fine-tune and optimize the DTC model's hyperparameters, particularly the DTFO and DTBW hybrid models. The primary aim was to enhance the predictive accuracy of these models. To evaluate the convergence of these optimization methods, two convergence curves (one related to G1 and the other related to G3) were utilized (see Fig. 3), tracking accuracy over 200 iterations. This curve visually demonstrated the evolution of Accuracy with each iteration, enabling the assessment of convergence progress and rate. In the case of G1 values, both models initially showed similar convergence rates of nearly 0.8, but the DTFO model ultimately achieved higher accuracy (almost 0.94). Notably, a linear pattern in the trend line around the 160-iteration mark indicated the optimal computational efficiency point for the DTFO model. On the other hand, regarding the G3 values, the DTFO model registered a lower convergence value at the beginning and a higher convergence value at the final iteration; it achieved a high convergence value of 0.92 at the final stage.
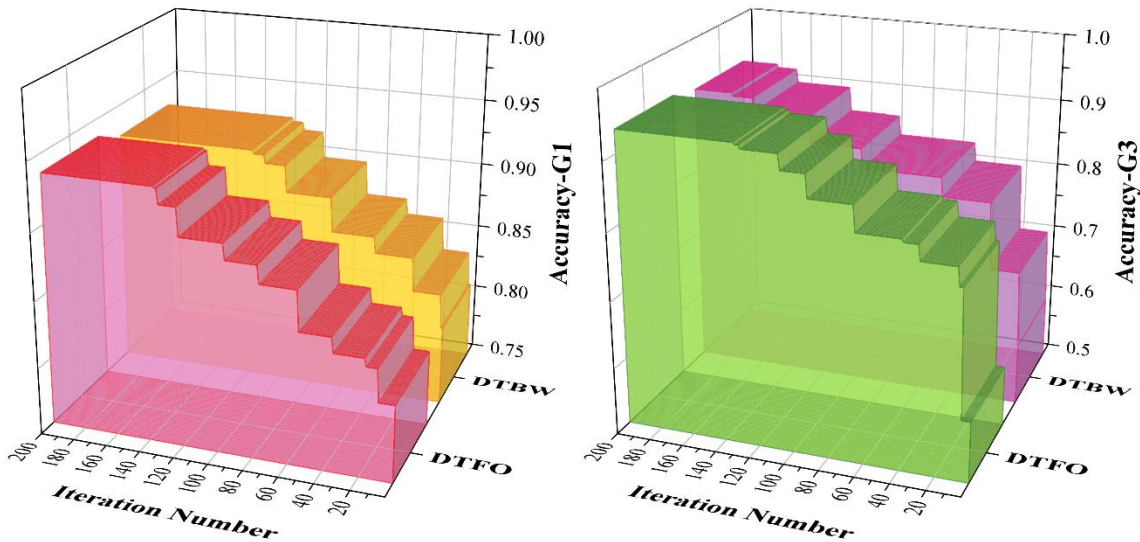


Fig. 3. Convergence of hybrid models.

## B. Hyperparameter

Table I displays the results of hyperparameter tuning for four different decision tree models, each associated with a specific target variable (G1 or G3). The hyperparameters include `max_depth` (maximum depth of the tree), `min_samples_split` (minimum samples required to split an internal node), `min_samples_leaf` (minimum samples required at a leaf node), and `max_leaf_nodes` (maximum number of leaf nodes). The values in each cell represent the chosen hyperparameter settings for the corresponding model and target variable. The hyperparameter tuning process aims to optimize the performance of the decision tree models in predicting student outcomes (G1 or G3).

The overall influence involves balancing model complexity and generalization. Higher values tend to lead to more complex models prone to overfitting, while lower values result in simpler models that generalize better. Hyperparameter tuning aims to find the optimal combination for the effective prediction of student outcomes.

TABLE I. RESULT OF HYPERPARAMETER

| Hyperparameter | Model (Target) | | | |
|---|---|---|---|---|
| | DTFO (G1) | DTBW (G1) | DTFO (G3) | DTBW (G3) |
| max_depth | 71 | 661 | 106 | 467 |
| min_samples_split | 0.001 | 0.209 | 0.001 | 0.116 |
| min_samples_leaf | 0.0005 | 0.0038 | 0.0005 | 0.0415 |
| max_leaf_nodes | 580 | 5 | 1270 | 4 |

## C. Comparing results of predictive models

This study focused on constructing three prediction models employing a classification approach to forecast students' exam performance in Mathematics and systematically improve their forthcoming grades. The models comprised a single Decision Tree Classification (DTC) and two optimized models using the Fox Optimization (FO) and the Black Widow Optimization (BWO). The dataset was split, allocating 70% for training and 30% for testing to assess their predictive performance. Table II and Fig. 4 illustrate the Accuracy, Precision, Recall, and F1-score for training, testing, and all phases across all models in predicting G1 and G3 scores.

- G1 Scores

Among the three models, the DTBW model exhibited superior training performance compared to the others, as evidenced by higher metric values during training than in the testing phase. The maximum metric values achieved by DTBW were 0.937 for all four metrics (Accuracy, Precision, Recall, and F1-Score). On the contrary, the DTC model obtained the lowest values, with 0.822 for Accuracy and Recall, 0.818 for Precision, and 0.82 for F1-Score.

- G3 Scores

Considering the mentioned models (DTC, DTFO, and DTBW), DTFO exhibited superior performance compared to the others, evident from its higher metric values. The maximum metric values achieved by DTFO were 0.934 for Accuracy and Recall and 0.935 for Precision and F1-Score. In contrast, the DTBW model obtained the lowest values, with 0.822 for Accuracy and Recall, 0.825 for Precision, and 0.823 for F1-Score.

TABLE II. RESULT OF PRESENTED MODELS

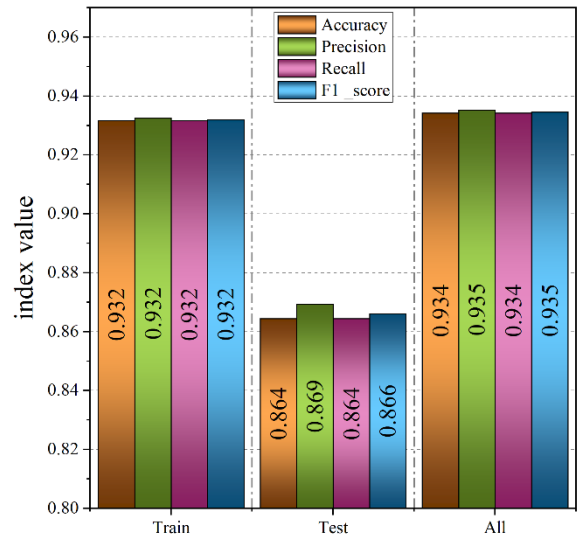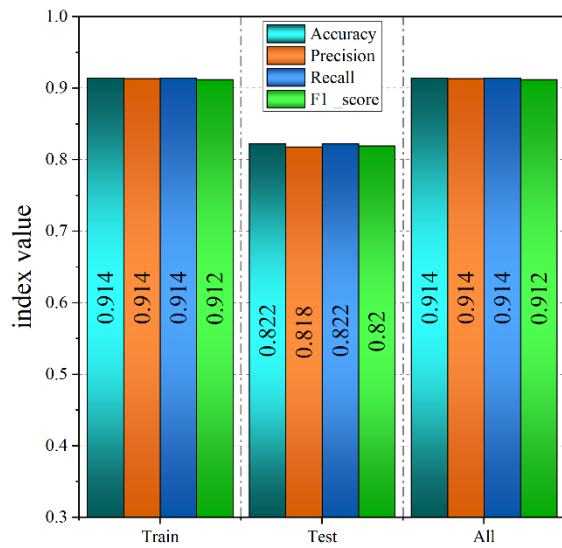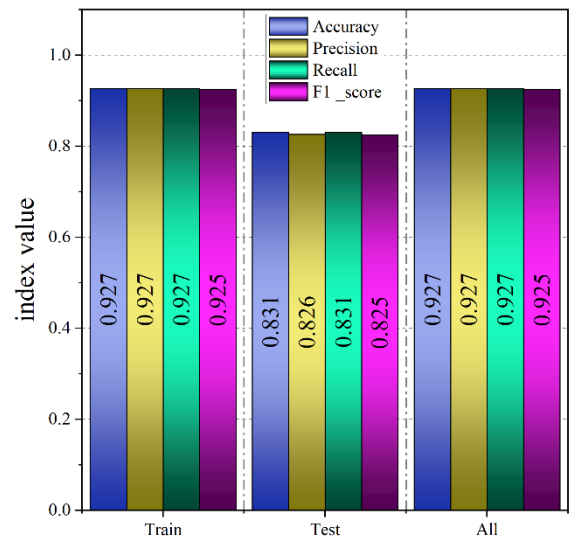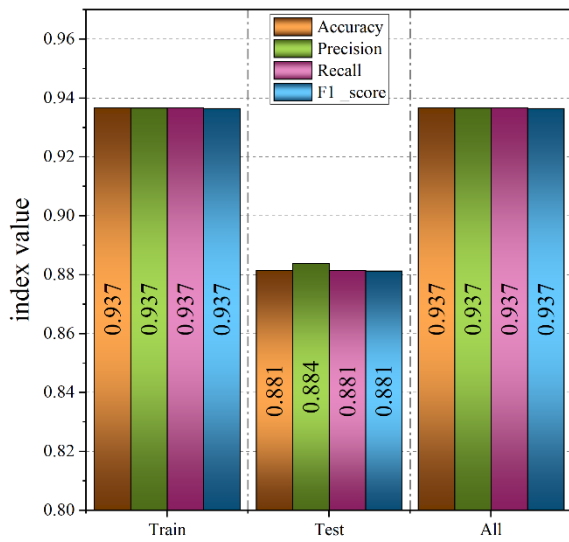| | Model | Phase | Index values | | | |
|---|---|---|---|---|---|---|
| | | | *Accuracy* | *Precision* | *Recall* | *F1 _core* |
| G1 | *DTC* | Train | 0.914 | 0.914 | 0.914 | 0.912 |
| | | Test | 0.822 | 0.818 | 0.822 | 0.820 |
| | | All | 0.914 | 0.914 | 0.914 | 0.912 |
| | *DTFO* | Train | 0.927 | 0.927 | 0.927 | 0.925 |
| | | Test | 0.831 | 0.826 | 0.831 | 0.825 |
| | | All | 0.927 | 0.927 | 0.927 | 0.925 |
| | *DTBW* | Train | 0.937 | 0.937 | 0.937 | 0.937 |
| | | Test | 0.881 | 0.884 | 0.881 | 0.881 |
| | | All | 0.937 | 0.937 | 0.937 | 0.937 |
| G3 | *DTC* | Train | 0.916 | 0.916 | 0.917 | 0.915 |
| | | Test | 0.856 | 0.854 | 0.856 | 0.852 |
| | | All | 0.916 | 0.916 | 0.917 | 0.917 |
| | *DTFO* | Train | 0.932 | 0.932 | 0.932 | 0.932 |
| | | Test | 0.864 | 0.869 | 0.864 | 0.866 |
| | | All | 0.934 | 0.935 | 0.934 | 0.935 |
| | *DTBW* | Train | 0.924 | 0.924 | 0.924 | 0.924 |
| | | Test | 0.822 | 0.825 | 0.822 | 0.823 |
| | | All | 0.924 | 0.924 | 0.924 | 0.924 |

Fig. 4. Column plot for the evaluation of developed models.

Following data processing and a comprehensive evaluation of the models' classification capabilities during the training and testing phases, 395 students were extensively examined based on their test results (G1 and G3 values). These students were categorized into four distinct groups: Poor (comprising students with scores ranging from 0 to 12), Acceptable (encompassing those with scores ranging from 12 to 14), Good (enrolling students with scores ranging from 14 to 20), and Excellent (comprising students with scores ranging from 16 to 20). The Index values for Precision, Recall, and F1-score are presented in Table III for G1 and Table IV for G3, which are used as evaluation metrics for assessing the classification performance of the developed models across the various student categories. A comparative analysis has been conducted in the subsequent section, considering each of these three Index values. As a result of this categorization, in the case of G1, 41 (10.38%) students were identified within the Excellent category, 54 (13.67%) within the Good category, 68 (17.21%) within the Acceptable category, and 232 (58.73%) within the Poor category. On the other hand, regarding G3 values, 40 (10.13%) students were identified within the Excellent category, 60 (15.19%) within the Good category, 62 (15.7%) within the Acceptable category, and 232 (58.73%) within the Poor category.

### D. Precision

- G1 Scores

The DTFO model demonstrated the highest values in the Good and Poor groups, achieving precision scores of 0.942 and 0.945, respectively. Conversely, the DTBW model obtained a maximum precision value of 0.947 for the Acceptable group. As for the excellent group, the DTC model outperformed others, attaining a precision score of 0.925.

- G3 Scores

The DTC model demonstrated the highest values in the Excellent and Acceptable categories, achieving precision scores of 0.922 and 0.898, respectively. On the other hand, the DTFO model obtained a maximum precision value of 0.974 for

the Poor group. As for the Good group, the DTBW model outperformed others, attaining a precision score of 0.9.

### E. Recall

- G1 Scores

The DTFO model displayed the highest scores in the Excellent, Good, and Acceptable groups, reaching 0.902, 0.907, and 0.897, respectively. When it comes to the Poor group, the DTBW model delivered the top performance with a recall score of 0.978.

- G3 Scores

In the Excellent and Good categories, the DTBW model demonstrated the highest values, achieving Recall values of 0.95 and 0.90, respectively. Furthermore, the DTC model obtained a maximum Recall value of 0.97 for the Poor group. While for the Acceptable group, the DTFO model outperformed others, attaining a score of 0.887.

### F. F1-score

- G1 Scores

A superior F1-score reflects the model's ability to balance precisely identifying positive cases (Precision) and encompassing all genuine positive cases (Recall). Upon considering all student categories, it becomes evident that the DTFO model demonstrated the highest values in the Good and Acceptable groups, achieving precision scores of 0.925 and 0.91, respectively. In addition, the DTBW model obtained a maximum F1-Score value of 0.956 for the Poor group. Finally, in the case of the Excellent group, the DTC model outperformed others, attaining an F1-Score of 0.914.

- G3 Scores

In the Excellent and Good categories, the DTBW model demonstrated the highest values, achieving F1-Score values of 0.927 and 0.90, respectively. Furthermore, the DTFO model outperformed others in the Poor category, attaining a score of 0.965. While for the Acceptable group, the DTC model obtained a maximum F1-Score value of 0.876.

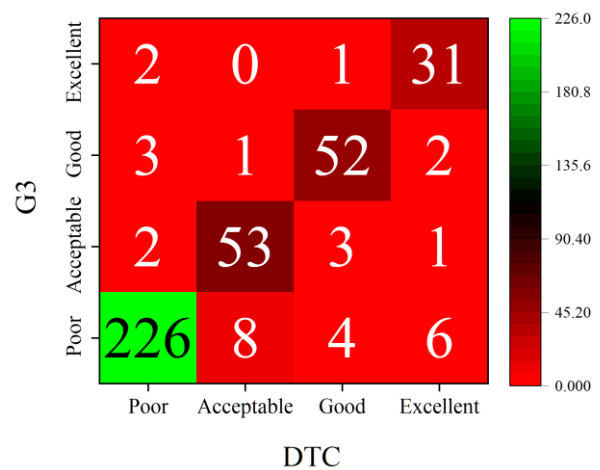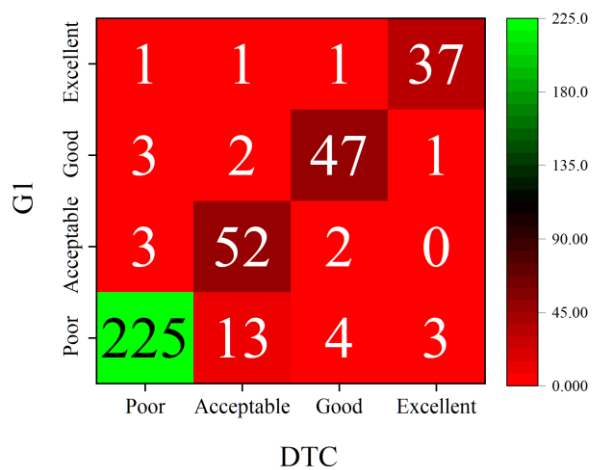TABLE III. EVALUATION INDEXES OF THE DEVELOPED MODELS' PERFORMANCE BASED ON G1

| Model | Grade | Index values | | |
|---|---|---|---|---|
| | | *Precision* | *Recall* | *F1-score* |
| **DTC** | *Excellent* | 0.925 | 0.902 | 0.914 |
| | *Good* | 0.887 | 0.870 | 0.879 |
| | *Acceptable* | 0.912 | 0.765 | 0.832 |
| | *Poor* | 0.918 | 0.970 | 0.943 |
| **DTFO** | *Excellent* | 0.902 | 0.902 | 0.902 |
| | *Good* | 0.942 | 0.907 | 0.925 |
| | *Acceptable* | 0.922 | 0.897 | 0.910 |
| | *Poor* | 0.945 | 0.961 | 0.953 |
| **DTBW** | *Excellent* | 0.881 | 0.902 | 0.892 |
| | *Good* | 0.906 | 0.889 | 0.897 |
| | *Acceptable* | 0.947 | 0.794 | 0.864 |
| | *Poor* | 0.934 | 0.978 | 0.956 |

TABLE IV.     EVALUATION INDEXES OF THE DEVELOPED MODELS' PERFORMANCE BASED ON G3

| Model | Grade | Index values | | |
|---|---|---|---|---|
| | | *Precision* | *Recall* | *F1-score* |
| DTC | *Excellent* | 0.912 | 0.775 | 0.838 |
| | *Good* | 0.897 | 0.867 | 0.881 |
| | *Acceptable* | 0.898 | 0.855 | 0.876 |
| | *Poor* | 0.926 | 0.970 | 0.948 |
| DTFO | *Excellent* | 0.884 | 0.950 | 0.916 |
| | *Good* | 0.898 | 0.883 | 0.891 |
| | *Acceptable* | 0.859 | 0.887 | 0.873 |
| | *Poor* | 0.974 | 0.957 | 0.965 |
| DTBW | *Excellent* | 0.905 | 0.950 | 0.927 |
| | *Good* | 0.900 | 0.900 | 0.900 |
| | *Acceptable* | 0.855 | 0.855 | 0.855 |
| | *Poor* | 0.952 | 0.944 | 0.948 |

The confusion matrix illustrated in Fig. 5 provides insights into accurately categorizing students into their respective grades and the misclassification into incorrect categories. In the case of G1 values, the DTFO model correctly categorized 37, 49, 61, and 223 students into Excellent, Good, Acceptable, and Poor classes, respectively, with only 25 students being misclassified. On the other hand, the DTBW and DTC models misclassified 29 and 34 students, respectively. Notably, misclassifications in the two optimized models primarily occurred between neighboring categories, such as 6 and 10 students for DTFO and DTBW, who were mistakenly placed in the Acceptable category instead of the Poor category. According to G3 values, the DTC model correctly categorized 31, 52, 53, and 223 students into Excellent, Good, Acceptable, and Poor classes, respectively, with 33 misclassified students. On the other hand, the DTBW and DTFO models misclassified 30 and 26 students, respectively. In the case of the single DTBW model, 9 students were inaccurately positioned in the Acceptable category instead of the Poor category.

The actual number of students falling into the Poor, Acceptable, Good, and Excellent categories was 232, 68, 54, and 41, respectively, for G1, while 233, 62, 60, and 40 for G3 values. Fig. 6 provides a visual representation of the student distribution across these categories based on measurement and classification model outcomes, facilitating a visual comparison. In the case of G1, the DTFO model exhibited the highest accuracy in correctly classifying students in the Acceptable, Good, and Excellent groups, identifying 61, 49, and 37 students accurately, respectively. In the case of the Poor category, the DTBW model outperformed the other models, correctly classifying 227 students. Regarding the G3 values, the DTFO model exhibited the highest accuracy in correctly classifying students into Acceptable and Excellent groups, identifying 55 and 38 students accurately. When considering the Poor category, the DTC model outperformed the other models, correctly classifying 226 students. Furthermore, according to the Good category, the DTBW model performed best, identifying 54 students correctly.
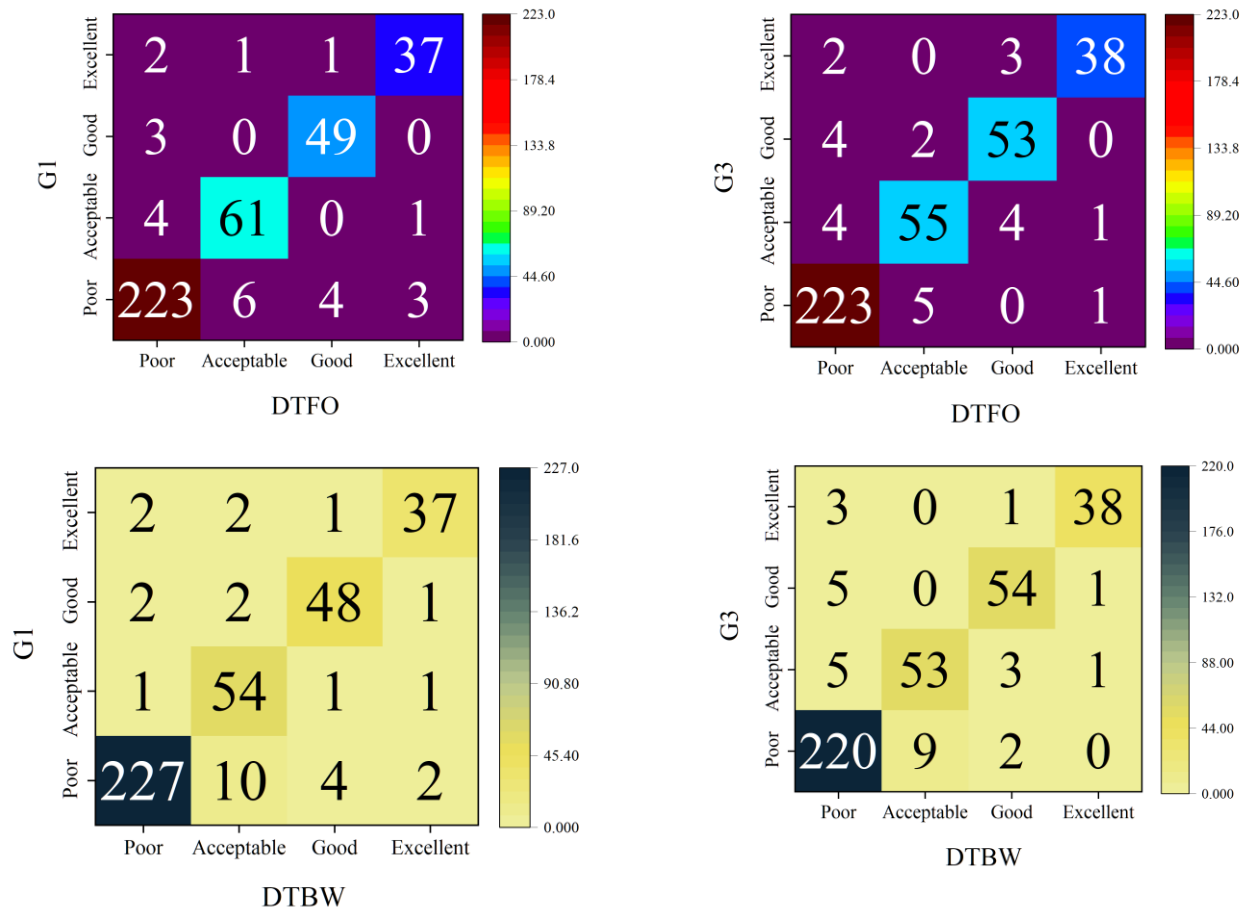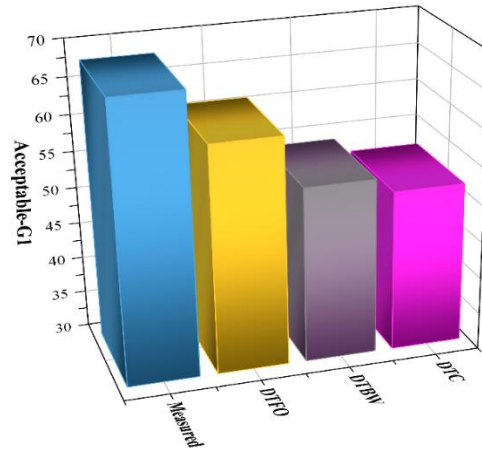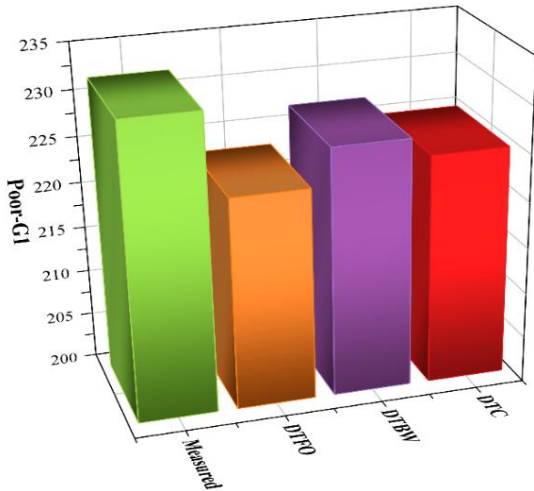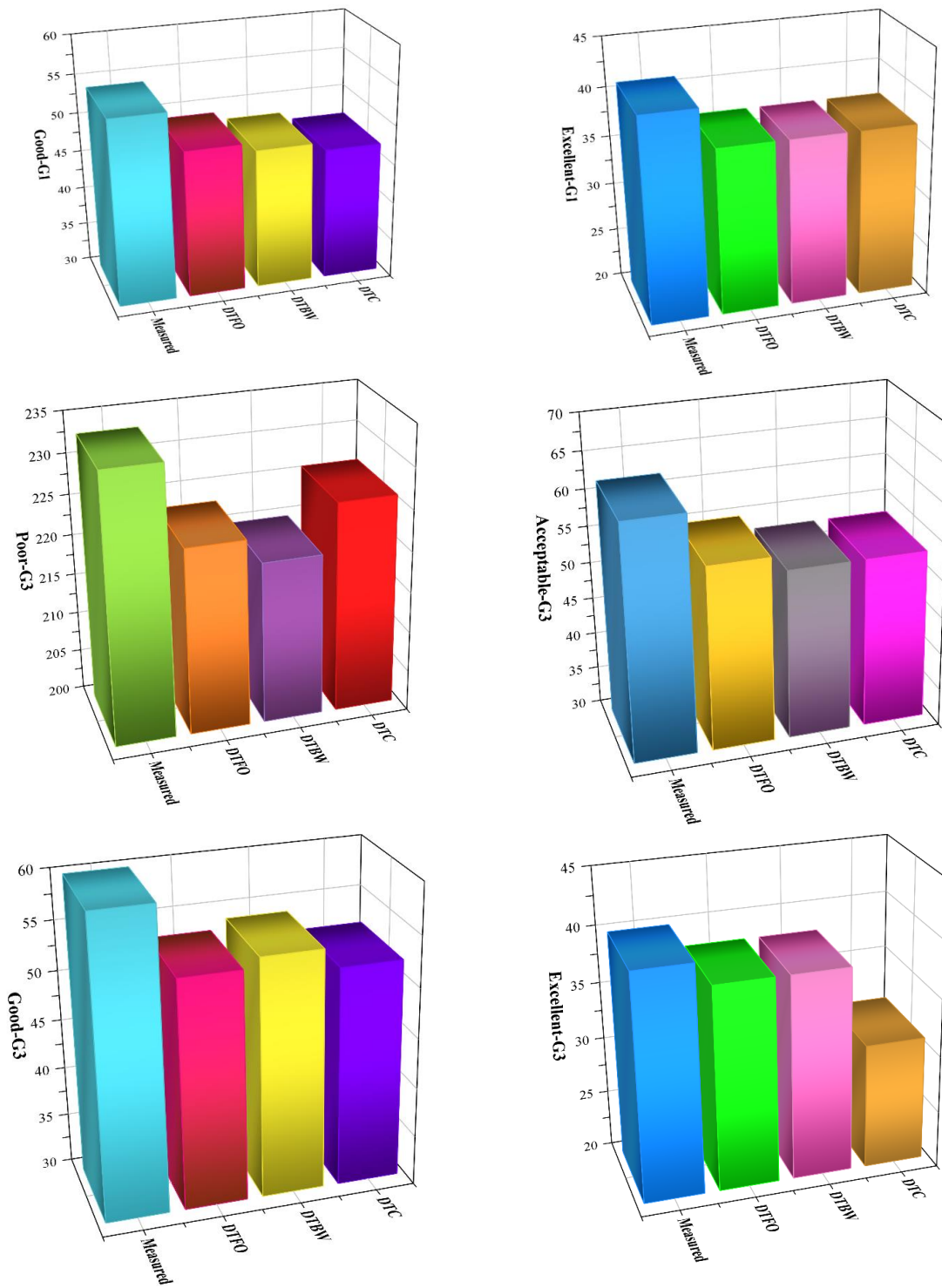
Fig. 5. Confusion matrix for each model's accuracy.

Fig. 6.    3D column plot for the developed models' accuracy compared to measured value.

## G. Sensitivity analyzes

SHAP (SHapley Additive exPlanations) values, derived from cooperative game theory, allocate feature contributions in ML models. They assess the impact of each feature on a model's prediction for a specific input, providing nuanced and interpretable insights. Adapted for use in ML, SHAP values offer a fair distribution of feature importance, aiding the interpretation of complex models by attributing predictions to individual features.

Fig. 7(a) reveals that "absences," "Freetime," "mother's job," and "Health" stand out as pivotal elements for anticipating G1 performance. Additionally, the plot highlights the fluctuating significance of these features across the four grade levels, indicating that the determinants influencing G1

scores are not uniform for all students. This underscores the variability in the impact of these factors across different grade levels.

On the other hand, for G3 in Fig. 7(b), it was observed that "absences," "Goout," "mother's education," and "mother's job " had the greatest impact on the model's output.
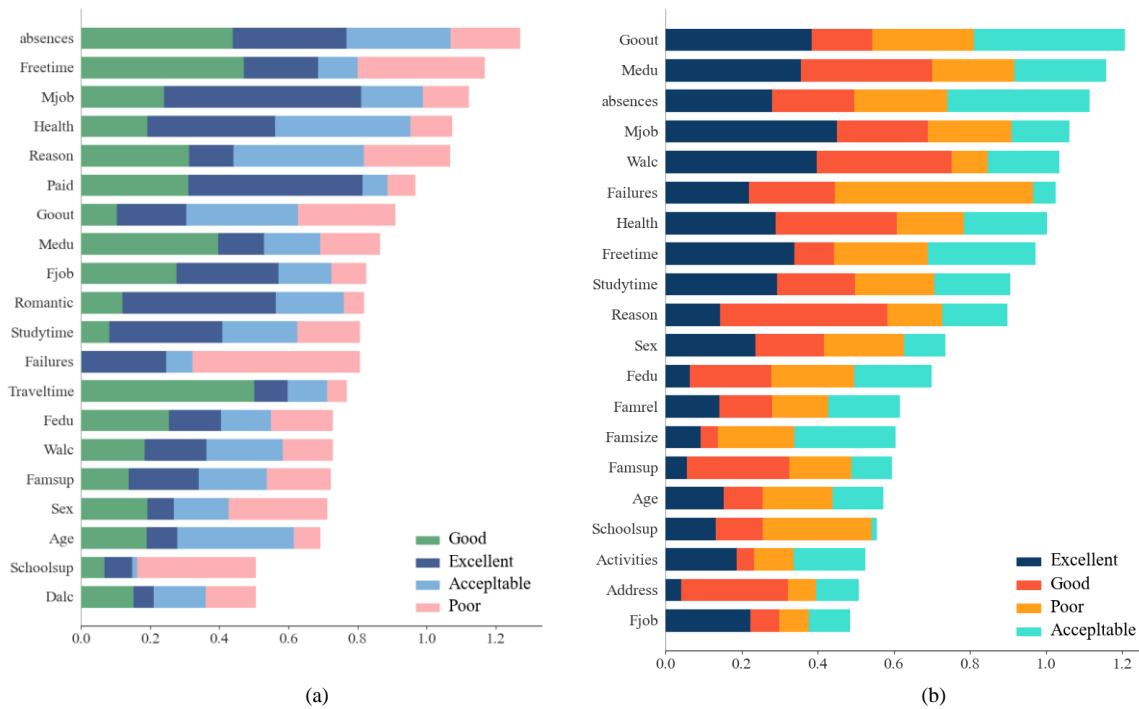


Fig. 7.    SHAP value for the impact of inputs on model's output a) G1 and b) G3.

## VI.    CONCLUSION

This research underscores the crucial significance of predictive models based on data in education. It stresses the need to consider qualitative and quantitative elements for predicting and evaluating students' academic performance. The findings offer valuable guidance for policymakers, educational institutions, and students, aiming to enhance future academic outcomes. The study demonstrates the effectiveness of data mining techniques such as clustering, classification, and regression in understanding and proactively tackling the diverse challenges encountered by undergraduate students. Furthermore, the research introduces an innovative approach by combining the Decision Tree Classification (DTC) model with optimization algorithms such as Fox Optimization (FO) and Black Widow Optimization (BWO). This advanced methodology illustrates how integrating machine learning techniques and optimization algorithms can elevate the Precision and effectiveness of predictive models. It provides a robust toolkit for addressing the evolving challenges in students' academic journeys. The study's thorough evaluation process, which included dividing the models into training and testing sets, reveals that these hybrid models have the potential to enhance the classification capabilities of the DTC model significantly. This enhancement is reflected in notable improvements in Accuracy and Precision. Upon analyzing the results, it has been observed that the potential to significantly enhance the classification capabilities of the DTC model by these hybrid models is increasingly recognized. Based on the results, it can be concluded that:

- In the case of G1 values, a marked improvement in Accuracy was achieved by applying FO and BWO optimization algorithms to the DTC model, with an increase of 1.42% and 2.51%, respectively. When the 395 students were categorized based on their final grades, the exceptional ability of the BWO to augment classification Accuracy became evident. Specifically, the DTBW model displayed an impressive Accuracy rate of 93.7%, accurately classifying the majority of students, whereas the DTFO and DTC models misclassified 6.33% and 8.6% of all students, respectively.

- With respect to G3 values, the improvement of Accuracy through the application of FO and BWO optimization algorithms to the DTC model was 1.96% for the application of FO and 0.87% for BWO. The DTFO model displayed an impressive Accuracy rate of 93.4%, accurately classifying the majority of students, whereas the DTBW and DTC models experienced misclassification rates of 7.59% and 8.35%, respectively.

The study sought to revolutionize academic performance prediction in education, assuming that predictive models significantly influence outcomes. Recognizing the holistic nature of student evaluation, it justified the importance of both qualitative and quantitative elements. Integration of machine learning with optimization algorithms was assumed to enhance predictive models, supported by literature. Standard practices in machine learning, such as thorough evaluation using training

and testing sets, were assumed to reflect model effectiveness. The assumption that misclassification rates indicate their direct measurement of prediction accuracy justified model performance. The study assumed that an increase in accuracy corresponded to improved classification capabilities, signifying enhanced predictions of students' final grades. Additionally, the assumption was made that optimization algorithms, specifically FO and BWO, led to marked improvements by fine-tuning decision tree models. Moreover, the research aimed to transform academic performance prediction in education, aligning with its overarching goal. Assumptions were strategically made to support this objective, including the significant influence of predictive models on academic outcomes, the importance of both qualitative and quantitative elements in predictions, and the enhancement of models through the integration of machine learning and optimization algorithms. Standard machine learning practices were assumed to reflect model effectiveness, with chosen metrics aligning with the study's goal of accurate predictions. The research assumed that improvements in accuracy corresponded to enhanced classification capabilities and that optimization algorithms led to marked improvements.

## REFERENCES

[1] S. Natek, M. Zwilling, Student data mining solution–knowledge management system related to higher education institutions, Expert Syst Appl 41 (2014) 6400–6407.

[2] Y. Zhao, C. Zhang, Y. Zhang, Z. Wang, J. Li, A review of data mining technologies in building energy systems: Load prediction, pattern identification, fault detection and diagnosis, Energy and Built Environment 1 (2020) 149–164.

[3] D. Kabakchieva, K. Stefanova, V. Kisimov, Analyzing university data for determining student profiles and predicting performance, in: Educational Data Mining 2011, 2010.

[4] C. Romero, S. Ventura, Educational data mining: A survey from 1995 to 2005, Expert Syst Appl 33 (2007) 135–146.

[5] C. Romero, S. Ventura, M. Pechenizkiy, R.Sj. Baker, Handbook of educational data mining, CRC press, 2010.

[6] R.S.J.D. Baker, K. Yacef, The state of educational data mining in 2009: A review and future visions, Journal of Educational Data Mining 1 (2009) 3–17.

[7] A. Ahmed, I.S. Elaraby, Data mining: A prediction for student's performance using classification method, World Journal of Computer Application and Technology 2 (2014) 43–47.

[8] E. Chandra, K. Nandhini, Knowledge mining from student data, European Journal of Scientific Research 47 (2010) 156–163.

[9] H.A.A. Hamza, P. Kommers, A review of educational data mining tools & techniques, International Journal of Educational Technology and Learning 3 (2018) 17–23.

[10] M.M.A. Tair, A.M. El-Halees, Mining educational data to improve students' performance: a case study, International Journal of Information 2 (2012).

[11] F. Ünal, Data mining for student performance prediction in education, Data Mining-Methods, Applications and Systems 28 (2020) 423–432.

[12] J.-M. Trujillo-Torres, H. Hossein-Mohand, M. Gómez-García, H. Hossein-Mohand, F.-J. Hinojo-Lucena, Estimating the academic performance of secondary education mathematics students: A gain lift predictive model, Mathematics 8 (2020) 2101.

[13] M.R. Apriyadi, D.P. Rini, Hyperparameter Optimization of Support Vector Regression Algorithm using Metaheuristic Algorithm for Student Performance Prediction, International Journal of Advanced Computer Science and Applications 14 (2023).

[14] C. Márquez-Vera, A. Cano, C. Romero, S. Ventura, Predicting student failure at school using genetic programming and different data mining approaches with high dimensional and imbalanced data, Applied Intelligence 38 (2013) 315–330.

[15] C. Romero, S. Ventura, Educational data mining: a review of the state of the art, IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews) 40 (2010) 601–618.

[16] B. Sekeroglu, K. Dimililer, K. Tuncal, Student performance prediction and classification using machine learning algorithms, in: Proceedings of the 2019 8th International Conference on Educational and Information Technology, 2019: pp. 7–11.

[17] A.K. Pal, S. Pal, Data mining techniques in EDM for predicting the performance of students, International Journal of Computer and Information Technology 2 (2013) 2279–2764.

[18] D. Kabakchieva, Student performance prediction by using data mining classification algorithms, International Journal of Computer Science and Management Research 1 (2012) 686–690.

[19] E. Osmanbegovic, M. Suljic, Data mining approach for predicting student performance, Economic Review: Journal of Economics and Business 10 (2012) 3–12.

[20] S. Huang, N. Fang, Predicting student academic performance in an engineering dynamics course: A comparison of four types of predictive mathematical models, Comput Educ 61 (2013) 133–145.

[21] L. Ramanathan, S. Dhanda, D.S. Kumar, Predicting students' performance using modified ID3 algorithm, International Journal of Engineering and Technology 5 (2013) 2491–2497.

[22] E.B. Costa, B. Fonseca, M.A. Santana, F.F. de Araújo, J. Rego, Evaluating the effectiveness of educational data mining techniques for early prediction of students' academic failure in introductory programming courses, Comput Human Behav 73 (2017) 247–256.

[23] F. Marbouti, H.A. Diefes-Dux, K. Madhavan, Models for early prediction of at-risk students in a course using standards-based grading, Comput Educ 103 (2016) 1–15.

[24] Y.-H. Hu, C.-L. Lo, S.-P. Shih, Developing early warning systems to predict students' online learning performance, Comput Human Behav 36 (2014) 469–478.

[25] D. Thammasiri, D. Delen, P. Meesad, N. Kasap, A critical assessment of imbalanced class distribution problem: The case of predicting freshmen student attrition, Expert Syst Appl 41 (2014) 321–330.

[26] M. Pandey, S. Taruna, A multi-level classification model pertaining to the student's academic performance prediction, Int J Adv Eng Technol 7 (2014) 1329.

[27] H. Sharma, S. Kumar, A survey on decision tree algorithms of classification in data mining, International Journal of Science and Research (IJSR) 5 (2016) 2094–2097.

[28] X. Luo, Efficient English text classification using selected machine learning techniques, Alexandria Engineering Journal 60 (2021) 3401–3409.

[29] D. Połap, M. Woźniak, Red fox optimization algorithm, Expert Syst Appl 166 (2021) 114107.

[30] V. Hayyolalam, A.A.P. Kazem, Black widow optimization algorithm: a novel meta-heuristic approach for solving engineering optimization problems, Eng Appl Artif Intell 87 (2020) 103249.

[31] J. Holland, Adaptation in natural and artificial systems, univ. of mich. press, Ann Arbor 7 (1975) 390–401.